

# Asymptotic distribution of $\Delta$ AUC, NRIs, and IDI based on theory of U-statistics

Olga V. Demler,<sup>a,\*†</sup>  Michael J. Pencina,<sup>b</sup>  Nancy R. Cook<sup>a</sup>  and Ralph B. D'Agostino Sr<sup>c</sup>

The change in area under the curve ( $\Delta$ AUC), the integrated discrimination improvement (IDI), and net reclassification index (NRI) are commonly used measures of risk prediction model performance. Some authors have reported good validity of associated methods of estimating their standard errors (SE) and construction of confidence intervals, whereas others have questioned their performance. To address these issues, we unite the  $\Delta$ AUC, IDI, and three versions of the NRI under the umbrella of the U-statistics family. We rigorously show that the asymptotic behavior of  $\Delta$ AUC, NRIs, and IDI fits the asymptotic distribution theory developed for U-statistics. We prove that the  $\Delta$ AUC, NRIs, and IDI are asymptotically normal, unless they compare nested models under the null hypothesis. In the latter case, asymptotic normality and existing SE estimates cannot be applied to  $\Delta$ AUC, NRIs, or IDI. In the former case, SE formulas proposed in the literature are equivalent to SE formulas obtained from U-statistics theory if we ignore adjustment for estimated parameters. We use Sukhatme–Randles–deWet condition to determine when adjustment for estimated parameters is necessary. We show that adjustment is not necessary for SEs of the  $\Delta$ AUC and two versions of the NRI when added predictor variables are significant and normally distributed. The SEs of the IDI and three-category NRI should always be adjusted for estimated parameters. These results allow us to define when existing formulas for SE estimates can be used and when resampling methods such as the bootstrap should be used instead when comparing nested models. We also use the U-statistic theory to develop a new SE estimate of  $\Delta$ AUC. Copyright © 2017 John Wiley & Sons, Ltd.

**Keywords:** AUC; NRI; IDI; risk prediction; U-statistics

## 1. An introduction and a motivating example

In current medical research, risk prediction is viewed as an objective way to assess the risk of a patient to develop a disease and is often used by clinicians in making treatment decisions. The Framingham [1] and ATP III models for 10-year risk of cardiovascular outcomes [2] and the Gail model for 5-year risk of breast cancer [3] are among the first widely used risk prediction models. Moreover, in recent years, risk-prediction models have played an increasingly important role in medical decision making and have been directly incorporated into updates of existing treatment guidelines. For instance, the U.S. Preventive Services Task Force recently issued updated guidelines on aspirin use in prevention of cardiovascular events [4]. Based on the results of a microsimulation model, that used the American College of Cardiology/American Heart Association (ACC/AHA) risk equations for 10-year cardiovascular disease risk [5]. Therefore, the quality of the performance of a risk prediction model is often crucial for assigning the most beneficial treatment and making correct policy decisions.

Risk prediction models are often evaluated in terms of calibration and discrimination. Discrimination measures how well a given model separates events from non-events; calibration measures the closeness of the model-based and observed risks of the outcome. The area under the receiver operating characteristics curve (AUC of ROC) [6,7] is a widely used measure of discrimination. In 2008, several new intuitively appealing measures of discrimination were introduced such as the net reclassification index (NRI) and integrated discrimination improvement (IDI) [8,9].

<sup>a</sup>Division of Preventive Medicine, Brigham and Women's Hospital, 900 Commonwealth Avenue, Boston, MA, 02115, U.S.A.

<sup>b</sup>Department of Biostatistics and Bioinformatics, Duke University, Durham, NC, 27708, U.S.A.

<sup>c</sup>Department of Mathematics and Statistics, Boston University, 111 Cummington Mall, Boston, MA, 02215, U.S.A.

\*Correspondence to: Olga V. Demler, Division of Preventive Medicine, Brigham and Women's Hospital, 900 Commonwealth Avenue, Boston, MA 02115, U.S.A.

†E-mail: olgademler@gmail.com

They rapidly gained popularity and at the time of writing this paper had been referenced more than 2800 times. Simple estimators for variance and asymptotic distributional behavior were proposed to allow construction of confidence intervals.

While some papers reported good validity of the methods for confidence intervals and variance estimators of  $\Delta\text{AUC}$ , NRIs, and IDI [8,10,11], others questioned their performance [10,12–14]. To illustrate these conflicting views, we ran some simulations and summarize the results in Table I. For two nested models with binary outcome and multivariate normal predictor variables, we compare observed and theoretical standard errors of  $\Delta\text{AUC}$ , three types of NRIs (continuous ( $\text{NRI}_{>0}$ ), 2-category NRI at event rate threshold ( $\text{NRI}(\text{r})$ ) and 3-category NRI (3cNRI)), and IDI. AUC is a measure of discrimination. It is equal to the probability that the risk of a randomly picked event is greater than for randomly picked non-event [6,7].  $\Delta\text{AUC}$  measures improvement in quality of discrimination between events and non-event by the new model relative to the old one [11].  $\text{NRI}_{>0}$ , another measure of discrimination, calculates the difference between fractions of correct and incorrect movements of predicted probabilities among events and adds to it a similar quantity calculated for non-events [9]. Categorical NRIs are similar to  $\text{NRI}_{>0}$  but consider only movements across categories.  $\text{NRI}(\text{r})$  uses two categories defined by event rate threshold [15]. 3cNRI uses three categories defined by any thresholds [16]. IDI combines average change in probabilities among events and among non-events [8]. For comparison, we included in Table I the regression coefficient ( $\beta$ ) for the new predictor variable  $x_2$ . The relative bias of standard error estimate is calculated as  $\frac{(\text{theoretical se} - \text{observed se})}{\text{observed se}} 100\%$ . Shaded areas in Table I indicate scenarios in which the relative bias is 5% or more in our simulations, while white areas indicate when standard errors have very low bias (<5%). Asymptotic theory developed for three of the five statistics performed very well in most situations, while the bias of the 3cNRI is comparable with that of the standard error estimator of the Kaplan–Meier survival probability (when sample size is small) [17], and the standard error estimator of the IDI has the strongest bias of the five statistics.

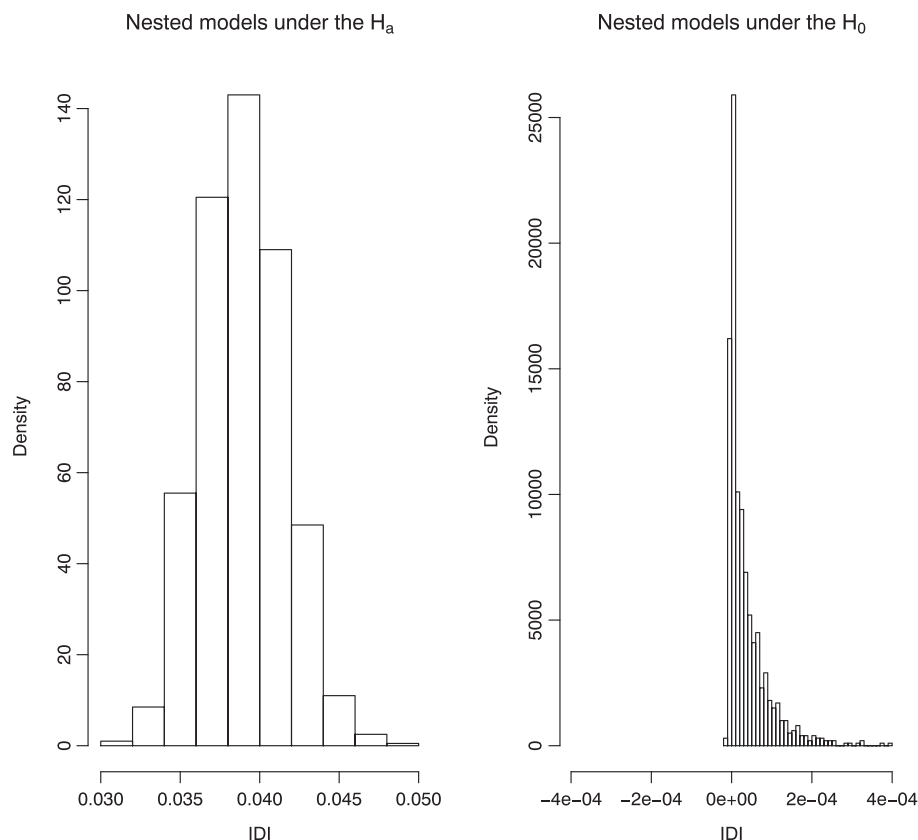
Confidence intervals for  $\Delta\text{AUC}$ ,  $\text{NRI}_{>0}$ , categorical NRIs, and IDI proposed to date rely on asymptotic normality [8,9,11,18,19]. In Figure 1, we show an example in which the IDI can be asymptotically normally distributed under the alternative hypothesis of meaningful effect (left panel) and right-skewed under the null hypothesis of no meaningful effect (right panel) [20].

This paper is a validity study of previously proposed asymptotic distribution results of  $\Delta\text{AUC}$ , IDI, and three types of NRIs (continuous ( $\text{NRI}_{>0}$ ), 2-category NRI at event rate threshold ( $\text{NRI}(\text{r})$ ), and 3cNRI) [8,9,11,18,19] when comparing two nested models. Using U-statistics theory, we explicitly specify conditions when asymptotic results are valid and when resampling methods such as the bootstrap should be used instead. These results help us disentangle several reports of the asymptotic distribution and performance of variance estimators of  $\Delta\text{AUC}$ , IDI, and three types of the NRI. The paper is structured as follows: Notation is introduced in Section 2; the main result is stated and proved in Section 3; in Sections 4 and 5, we apply theoretical findings to the Framingham Heart Study (FHS) data; and the implications of these findings are discussed in Section 6.

**Table I.** Relative bias (%) of standard errors of  $\Delta\text{AUC}$ ,  $\text{NRI}_{>0}$ , 3-category NRI, and IDI.

Effect size	Size	Relative bias (%) of the standard error estimate of					
		$\beta_{x_2}$	$\Delta\text{AUC}$	$\text{NRI}_{>0}$	$\text{NRI}(\text{r})$	3cNRI	IDI
0	30,000	3	7	29	−7	−7	−31
0	2,000	0	8	26	−12	−13	−37
0	500	0	8	26	−12	−13	−37
0.2	30,000	3	1	−3	1	−5	−42
0.2	2,000	1	−1	−3	−1	−10	−45
0.2	500	1	−1	−3	−1	−10	−45
0.7	30,000	2	0	−1	0	−26	−38
0.7	2,000	1	3	−2	1	−26	−39
0.7	500	1	3	−2	1	−26	−39

( $\text{rel.bias} = \frac{\text{SE}_{\text{formula-based}} - \text{SE}_{\text{bootstrap}}}{\text{SE}_{\text{bootstrap}}} \times 100\%$ ). We evaluated the performance of two nested risk prediction models: a logistic regression model with two multivariate normally distributed predictor variables ( $x_1$  and  $x_2$ ) and a baseline logistic regression model with only one of the predictors ( $x_1$ ). We considered several simulations scenarios: effect size by the new predictor ( $x_2$ ) of 0, 0.2, and 0.7; effect size by commonly used predictor variable ( $x_1$ ) is 0.7; sample sizes of 30,000, 2,000, and 500 observations; 0.1 event rate; and  $B = 1,000$  simulated datasets. 2% and 10% cutoffs were used for 3cNRI calculation.



**Figure 1.** Histograms of IDI when comparing nested models under the alternative (left panel) and under the null (right panel).  $x_1, x_2$  are predictors from the full model;  $x_1$  is the predictor from the reduced model. Left panel: simulated nested models under the alternative  $x_1, x_2 \text{ ID} = 1 \sim N(\mu, \Sigma)$  and  $x_1, x_2 \text{ ID} = 0 \sim N(0, \Sigma)$ . Right panel: simulated nested models under the null  $x_1 \text{ ID} = 1 \sim N(\mu, \sigma^2)$ ,  $x_1 \text{ ID} = 0 \sim N(0, \sigma^2)$ , and  $x_2 \sim N(0, \sigma^2)$ .  $x_2$  is an uninformative predictor.

## 2. Notation

Let  $D$  be an outcome of interest, with  $D = 1$  for events and  $D = 0$  for non-events. Our goal is to predict the event status using  $p$  predictor variables. Conditioning on the event status, predictor variables follow two (potentially different) distribution functions:  $\mathbf{x}|D = 0 \sim \mathbf{F}(\cdot)$ ,  $\mathbf{y}|D = 1 \sim \mathbf{G}(\cdot)$ . Assume that for each of  $N$  patients, their disease status  $D$  and vector of predictor variables are available. There are  $n_0$  non-events and  $n_1$  events. The prediction based on the full set of  $p$  predictor variables is to be compared with that based on a reduced number of predictor variables,  $p - 1$ . We assume that the linear model is true and that one of the linear models for binary outcome is employed (logistic regression, linear discriminant analysis (LDA), etc). We use this model to estimate linear coefficients in order to combine multiple predictor variables into one metric, the risk score. Unless otherwise specified, we assume that the models are nested, so the new model adds  $k$  new predictors to the old model. The regression technique of choice produces coefficients estimates  $\mathbf{a}^{*'} = (a_1^*, \dots, a_{p-k}^*, 0, \dots, 0)$  (reduced model) and  $\mathbf{a}' = (a_1, \dots, a_p)$  (full model). Corresponding risk scores are calculated as  $\mathbf{a}'\mathbf{x}$  and  $\mathbf{a}^{*'}\mathbf{x}$  for non-events and  $\mathbf{a}'\mathbf{y}$  and  $\mathbf{a}^{*'}\mathbf{y}$  for events, with the symbol  $*$  always denoting the reduced model. We sought to test whether the risk prediction model with  $p$  predictors performs better than the model with only the first  $p-k$  predictors. We will consider  $\Delta\text{AUC}$ , three varieties of the NRI, and the IDI as measures of model performance. They are often used in current medical research on risk prediction. Analysis of their performance, advantages, and disadvantages is an active area of methodological research on risk prediction. In the following, we review standard formulas [8,9,11] for  $\Delta\text{AUC}$ , continuous NRI ( $\text{NRI}_{>0}$ ), 3cNRI, NRI(r), and IDI.

### 2.1. $\Delta\text{AUC}$

The area under ROC curve (AUC) can be interpreted as the probability that the risk score of a randomly picked event is higher than a randomly picked non-event. The AUC is estimated by the Mann–Whitney

statistic [6,7] – a non-parametric unbiased estimator, often referred to as the *c*-statistic [21,22] and can be written as follows:

$$AUC = \frac{1}{n_0 n_1} \sum_{i=1}^{n_0} \sum_{j=1}^{n_1} I[a'x_i < a'y_j], \text{ where } I[\cdot] \text{ is the indicator function.}$$

The AUC for the reduced model is as follows:  $AUC^* = \frac{1}{n_0 n_1} \sum_{i=1}^{n_0} \sum_{j=1}^{n_1} I[a^*x_i < a^*y_j]$ .

Then  $\Delta AUC$  is as follows:

$$\Delta AUC = \frac{1}{n_0 n_1} \sum_{i=1}^{n_0} \sum_{j=1}^{n_1} I[a'x_i < a'y_j] - \frac{1}{n_0 n_1} \sum_{i=1}^{n_0} \sum_{j=1}^{n_1} I[a^*x_i < a^*y_j]$$

$\Delta AUC$  is one of the most widely used measures of discrimination.

## 2.2. Continuous NRI ( $NRI_{>0}$ )

$NRI_{>0}$  [9] is the difference of proportions of individuals with events and non-events whose predicted probabilities moved up:

$$\begin{aligned} NRI_{>0} &= \frac{\sum_{i=1}^{n_1} \text{Sign}[p_{\text{new ev}} - p_{\text{old ev}}]}{n_1} - \frac{\sum_{i=1}^{n_0} \text{Sign}[p_{\text{new nonev}} - p_{\text{old nonev}}]}{n_0} \\ &= \frac{\# \text{events up}}{n_1} - \frac{\# \text{nonevents up}}{n_0} \end{aligned}$$

## 2.3. Three-category NRI

Three-category NRI [16] is very close to the original definition of categorical NRI [9] but takes into account the size of the jump from category to category (number of categories moved). It is defined as follows:

$$3cNRI = \frac{1}{n_1} \sum_{i=1}^{n_1} \# \text{categories up}_i - \# \text{categories down}_i - \frac{1}{n_0} \sum_{j=1}^{n_0} \# \text{categories up}_j - \# \text{categories down}_j$$

This definition of categorical NRI is preferable over its original 2008 version [8], because of several attractive properties [16], including the fact that  $3cNRI = 0$  if marginal cells of the reclassification table stay the same for the two models. By using weights, it treats jumps across one versus two categories differently, and the event rate has a limited impact on the magnitude of the  $3cNRI$ . Therefore, it successfully resolves several criticisms of the original definition of categorical NRI [23,24].

## 2.4. NRI at the event rate ( $NRI(r)$ )

In their 2016 paper, Pencina, Steyerberg, and D'Agostino [15] investigate the properties of a 2-category NRI with categories defined by the proportion of cases in the sample (*r*) and show that it has several advantages: like  $\Delta AUC$ , it is invariant to the event rate and has intuitive interpretation as the proportion of correct reclassifications.

## 2.5. IDI

IDI [8] is defined as

$$IDI = \frac{\sum_{i=1}^{n_1} p_{\text{new ev } i} - p_{\text{old ev } i}}{n_1} - \frac{\sum_{i=1}^{n_0} p_{\text{new nonev } i} - p_{\text{old nonev } i}}{n_0}$$

Discrimination improvement is related asymptotically to the rescaled Brier score and to the difference in discrimination slope [25]. We mentioned some criticisms of IDI earlier, and in the following, we address some of them.

Now we can formulate the following null hypotheses for the six statistics defined earlier:

$$\begin{aligned} H_0^{AUC} : \Delta AUC = 0 \quad &\text{vs} \quad H_a^{AUC} : \Delta AUC \neq 0 \\ H_0^{NRI} : NRI_{>0} = 0 \quad &\text{vs} \quad H_a^{NRI} : NRI_{>0} \neq 0 \\ H_0^{NRI} : NRI(r) = 0 \quad &\text{vs} \quad H_a^{NRI} : NRI(r) \neq 0 \\ H_0^{3cNRI} : 3cNRI = 0 \quad &\text{vs} \quad H_a^{3cNRI} : 3cNRI \neq 0 \\ H_0^{IDI} : IDI = 0 \quad &\text{vs} \quad H_a^{IDI} : IDI \neq 0 \end{aligned} \quad (1)$$

Pepe *et al.* [26,27] showed that each of the five hypotheses in (1) is equivalent to testing the significance of the set of the new predictors in the new regression model (2).

$$H_0 : a_{p-k+1}, \dots, a_p = 0 \quad \text{vs} \quad H_a : a_{p-k+1}, \dots, a_p \neq 0 \quad (2)$$

Therefore, when we consider data under the null, we can without loss of generality assume that the null is formulated in terms of non-significance of the linear coefficient by the new predictor variable, that is, the hypothesis in (2).

### 3. Main result

We formulate our main results as follows.

$\Delta AUC$ ,  $NRI_{>0}$ ,  $NRI(r)$ ,  $3cNRI$ , and  $IDI$ :

STATEMENT 1. *are generalized U-statistics with estimated parameters.*

STATEMENT 2. *belong to non-degenerate subclass if and only if they compare any non-nested models or nested models under the alternative hypothesis in (2). As non-degenerate U-statistics,*

- they follow normal distribution asymptotically.*
- Available variance formulas are algebraically equal to the variance estimators provided by U-statistics theory if we ignore adjustment for estimated parameters.*
- Variance of  $\Delta AUC$ ,  $NRI_{>0}$ , and  $NRI(r)$  does not need to be adjusted for estimated parameters if predictor variables are normally distributed.*
- Variance of  $IDI$  and 3-category  $NRI$  should always be adjusted for estimated parameters.*

STATEMENT 3. *belong to the **degenerate subclass** if and only if they compare nested models under the null hypothesis in (2). As degenerate U-statistics, they do not follow normal distribution and available variance estimators do not apply for them.*

#### 3.1. $\Delta AUC$ , $NRI_{>0}$ , $NRI(r)$ , $3cNRI$ , and $IDI$ belong to the U-statistics family

In Appendix B, we prove Statement 1 showing that statistics considered in this paper belong to a U-statistics family [28]. Rigorous asymptotic distribution theory of U-statistics has been developed by Hoeffding [29], Lehman [30], Sukhatme [31], and others. The form of the U-statistics' distribution depends on whether the U-statistics are degenerate. Non-degenerate U-statistics are normally distributed, and formulas for their standard errors are available. Degenerate U-statistics are distributed as an infinite sum of weighted Chi-square random variables, and derivation of their standard error is challenging.

In Appendix B, we show that  $\Delta AUC$ ,  $NRI_{>0}$ ,  $NRI(r)$ ,  $3cNRI$ , and  $IDI$  are degenerate if and only if they compare nested models under the null. In all other situations, they belong to the non-degenerate class of U-statistics. Degeneracy and non-degeneracy conditions are listed in Table II.

Degenerate and non-degenerate U-statistics form very different classes in terms of their asymptotic behavior. In the following sections, we will consider these two situations separately.

#### 3.2. Non-degenerate case

$\Delta AUC$ ,  $NRI$ s, and  $IDI$  are non-degenerate if they evaluate the performance of two non-nested models or of nested models under the alternative. This is the most practically interesting case because only in this

**Table II.** Non-degeneracy conditions of  $\Delta\text{AUC}$ ,  $\text{NRI}_{>0}$ ,  $\text{NRI}(r)$ ,  $3\text{cNRI}$ , and  $\text{IDI}$ .

	Models are under the null	Models are under the alternative
Nested models	Degenerate*	Non-degenerate
Non-nested models	Always non-degenerate	

\*Null is defined as in (2) in the previous section.  $H_0: a_p - k + 1, \dots, a_p = 0$ .

situation we need to construct confidence intervals for  $\Delta\text{AUC}$ , NRIs, and IDI. Hoeffding [29] and Lehman [30] showed that non-degenerate U-statistics are asymptotically normally distributed. U-statistics theory also provides their variance formulas [28] but notes that variances should be adjusted for estimated parameters. Adjustment has been studied by Sukhatme [32], Randles [33], and de Wet [34] and is summarized in [28].

### 3.3. Available variance estimators are identical to U-statistics theory-based variance estimators if we ignore an adjustment for estimated parameters

In the Appendix, we derived variances of  $\Delta\text{AUC}$ , NRIs, and IDI based on U-statistics theory, ignoring adjustment for estimated parameters and presented them in Table III. The standard errors of  $\text{NRI}_{>0}$  and  $\text{NRI}(r)$  based on the U-statistics theory are exactly the same as the ones derived by Pencina *et al.* in [10,11]. The standard error formula for  $\Delta\text{AUC}$  is new. It is equal to the variance of the change in ranks. This representation is more intuitive, but it assumes no tied ranks.

U-statistics theory adds one more layer to variance calculations, namely that when U-statistic relies on estimated parameters, its variance in general should be adjusted for estimated parameters. In many cases  $\Delta\text{AUC}$ , NRIs, and IDI rely on estimated parameters (linear coefficients of regression models), their variances may need to be adjusted for estimated parameters, or we need to show that such adjustment is not necessary. In the following section, we prove that for some of the statistics under certain assumptions, adjustment for estimated parameters is unnecessary.

### 3.4. Variances of $\Delta\text{AUC}$ , $\text{NRI}_{>0}$ , and $\text{NRI}(r)$ do not need to be adjusted for estimated parameters if predictor variables are normally distributed

Sometimes, adjustment for estimated parameters can be avoided. Sukhatme [32], Randles [33], and de Wet [34] showed that adjustment for estimated parameters is unnecessary if and only if a certain condition is met [28]. In the following, we check this condition and show that under normality of predictor variables, standard error estimates of  $\Delta\text{AUC}$ , continuous NRI, and  $\text{NRI}(r)$  do not need to be adjusted for estimated parameters.

**Table III.** Variance formulas in non-degenerate case, unadjusted for estimated parameters.

	$\hat{\sigma}^2$ , ignoring the adjustment for estimated parameters	Requires adjustment?
$\hat{\sigma}_{\Delta\text{AUC}}^2$ no tied ranks	$\frac{\text{Var}(\text{rank}_e^*(a^T x_i) - \text{rank}_e(a^T x_i))}{n_0} + \frac{\text{Var}(\text{rank}_{ne}^*(a^T y_j) - \text{rank}_{ne}(a^T y_j))}{n_1}$	No
$\hat{\sigma}_{\Delta\text{AUC}}^2$ tied ranks	Use DeLong formula [11]	No
$\hat{\sigma}_{\text{NRI}_{>0}}^2$	$\frac{\hat{p}_{ne}^{up}(1 - \hat{p}_{ne}^{up})}{n_0} + \frac{\hat{p}_e^{up}(1 - \hat{p}_e^{up})}{n_1}$	No
$\hat{\sigma}_{\text{NRI}(r)}^2$	$\frac{\hat{p}_{ne}^{up} + \hat{p}_{ne}^{down} - (\hat{p}_{ne}^{up} - \hat{p}_{ne}^{down})^2}{n_0} + \frac{\hat{p}_{ev}^{up} + \hat{p}_{ev}^{down} - (\hat{p}_{ev}^{up} - \hat{p}_{ev}^{down})^2}{n_1}$	No
$\hat{\sigma}_{3\text{cNRI}}^2$	$\frac{4(\hat{p}_{ne}^{2up} + \hat{p}_{ne}^{2down}) + \hat{p}_{ne}^{1up} + \hat{p}_{ne}^{1down} - (2(\hat{p}_{ne}^{2up} - \hat{p}_{ne}^{2down}) + \hat{p}_{ne}^{1up} - \hat{p}_{ne}^{1down})^2}{n_0} + \frac{4(\hat{p}_{ev}^{2up} + \hat{p}_{ev}^{2down}) + \hat{p}_{ev}^{1up} + \hat{p}_{ev}^{1down} - (2(\hat{p}_{ev}^{2up} - \hat{p}_{ev}^{2down}) + \hat{p}_{ev}^{1up} - \hat{p}_{ev}^{1down})^2}{n_1}$	Yes
$\hat{\sigma}_{\text{IDI}}^2$	$\frac{\text{Var}(\Delta\text{predp}(x_i))}{n_0} + \frac{\text{Var}(\Delta\text{predp}(y_j))}{n_1}$	Yes



**STATEMENT 2.C** If  $\Delta AUC$ ,  $NRI_{>0}$ , and  $NRI(r)$  when comparing nested models are non-degenerate (Table II) and if predictor variables are normally distributed, then standard errors of  $\Delta AUC$ , continuous  $NRI$ , and  $NRI(r)$  do not need to be adjusted for estimated parameters.

*Proof*

We restate here the condition for adjustment for estimated parameters:

*Sukhatme–Randles–de Wet condition:*

*Standard errors for a U-statistic with estimated parameters do not need to be adjusted for estimated parameters if and only if the derivative of the expected value of the U-statistic with respect to parameters is zero.*

For example, for  $\Delta AUC$ , this condition is written as  $\frac{\partial}{\partial a} E[\Delta AUC] = 0$ .

In our assumptions, predictors are normally distributed; therefore, LDA is the most efficient way to estimate regression coefficients [35]. Su and Liu [36] also showed that under these assumptions, LDA coefficients maximize Mahalanobis distance [37] ( $M^2$ ) between risk scores of events and non-events. Therefore, the gradient of Mahalanobis distance with respect to parameters is zero. When comparing nested models,  $\Delta AUC$  is a function of the Mahalanobis distances ( $\Delta AUC = \Phi\left(\sqrt{\frac{M_p^2}{2}}\right) - \Phi\left(\sqrt{\frac{M_{p-k}^2}{2}}\right)$ ) [36], where  $p$  is the number of predictor variables in a model. Hence, the gradient of  $\Delta AUC$  with respect to parameters is zero as well. Therefore, the standard error of  $\Delta AUC$  under the assumption of normality of predictor variables does not need to be adjusted for estimated parameters.

Similarly, we can use a closed-form formula for  $NRI_{>0}$  [38] for nested models: ( $NRI_{>0} = 4\Phi\left(\frac{\sqrt{M_p^2 - M_{p-k}^2}}{2}\right) - 2$ ) to show that gradient of  $NRI_{>0}$  is also zero at the LDA coefficients. Therefore, the standard error of  $NRI_{>0}$  also does not need to be adjusted for estimated parameters.

Pencina, Steyerberg, and D'Agostino [15] showed that  $NRI(r)$  under normality assumptions when comparing nested models can be written as follows:

$$NRI(r) = 2 \cdot \left( \Phi\left(\frac{\sqrt{M_p^2}}{2}\right) - \Phi\left(\frac{\sqrt{M_{p-k}^2}}{2}\right) \right).$$
 The same reasoning can be applied to  $NRI(r)$  to show that  $\frac{\partial}{\partial a} NRI(r) = 0$ . Therefore,  $NRI(r)$  does not need to be adjusted for estimated parameters under the assumptions of this statement.

Q.E.D.

**STATEMENT 2.D** Variances of IDI and 3cNRI should always be adjusted for estimated parameters.

Note that the IDI and 3cNRI also can be expressed in closed form under normality of predictor variables [16,38] (please see the Appendix), but their closed-form expression does not rely exclusively on the Mahalanobis distance. It also depends on the estimated rate of events, which becomes one of the parameters. Under normality of predictor variables, LDA solution maximizes Mahalanobis distance, and therefore, the derivative of  $M_p^2$  with respect to regression parameters is zero. However, there is no such result for partial derivative of the closed-form formulas of IDI and 3cNRI with respect to event rate. Derivatives of closed-form formulas of 3cNRI and IDI with respect to event rate were calculated in Appendix D. Both derivatives are nonlinear in  $r$ , and both are in general non-zero. For example, derivatives of 3cNRI and IDI are 2.02 and 1.04 correspondingly for event rate observed in FHS of 7.65%, when comparing models with Mahalanobis distances of 0.7 and 0.8 and using 5% and 7.5% cutoffs to calculate 3cNRI. Therefore, the Sukhatme–Randles–deWet condition is not satisfied for IDI and 3cNRI, and standard errors of IDI and 3cNRI should be adjusted for estimated parameters.

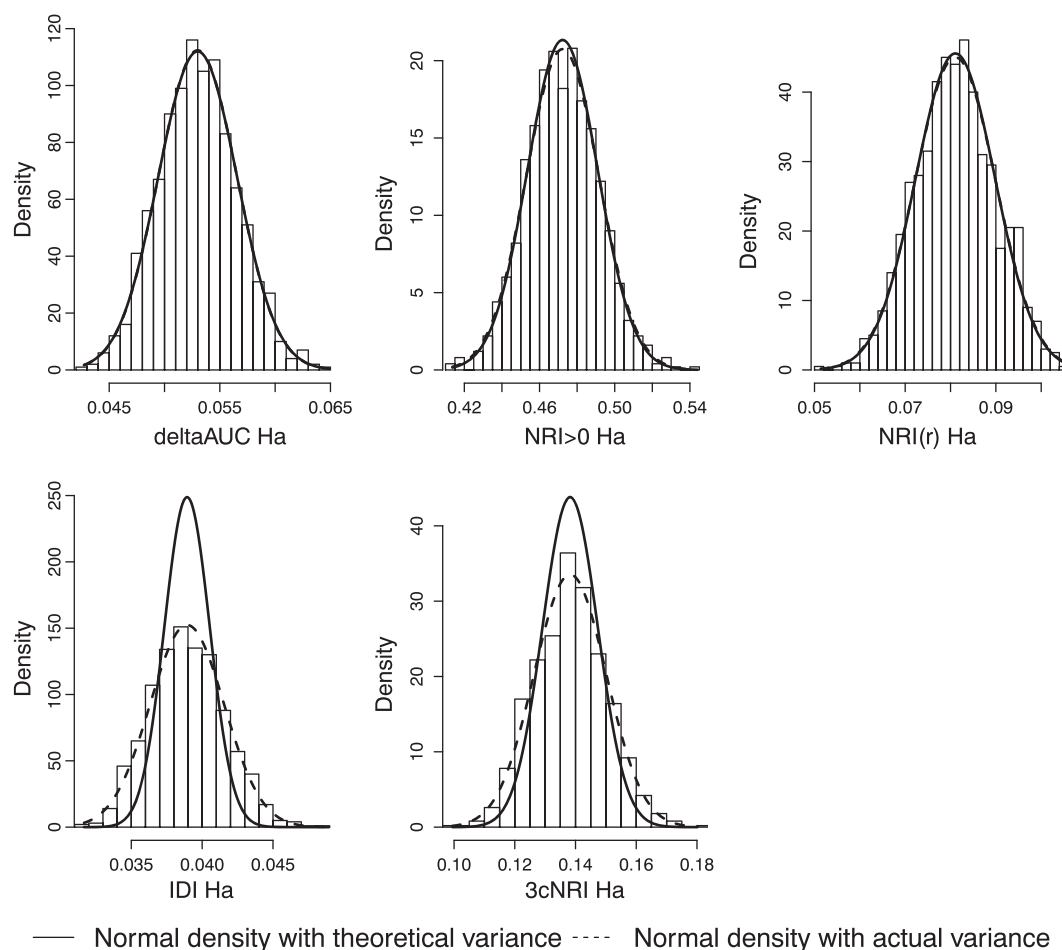
Our empirical results in Table I support the main theoretical results proven in this paper. Variances of  $\Delta AUC$ ,  $NRI_{>0}$ , and  $NRI(r)$  calculated from unadjusted formulas have on average very small relative bias compared with those of the IDI and 3cNRI whose variances must be adjusted for estimated parameters.

Also, the top three rows of Table I are calculated for the degenerate case (when comparing two nested models under the null). All five statistics are degenerate, and theoretical formulas for their variance estimator are not applicable: Existing variance formulas have strong bias for all five statistics when comparing nested models under the null.

To illustrate further the main theoretical findings of this paper, we simulated a binary model with normally distributed predictor variables. In Figure 2, we plot histograms of  $\Delta\text{AUC}$ , NRIs, and IDI calculated for nested models under the alternative and overlay two normal distribution curves with empirical (dotted line) and theoretical (solid) variances.  $\Delta\text{AUC}$ ,  $\text{NRI}_{>0}$ , and  $\text{NRI}(r)$  are in the top row. They do not need to be adjusted for estimated parameters, and the dotted and solid curves almost completely overlap. IDI and  $3\text{cNRI}$  are in the bottom row. They require adjustment for estimated parameters, and the two curves do not overlap because the theoretical variance is an incorrect estimate of the actual variance of  $3\text{cNRI}$  and IDI.

### 3.5. Statement 2.C and 2.D for logistic regression and non-normal data

We showed in the proof of Statement 2.C that by estimating parameters with LDA, we ensured that Sukhatme–Randles–deWet condition holds true. What would happen if we had used logistic regression to estimate parameters instead of the LDA? To use theoretical variance formulas, we need to show that adjustment for parameters estimated by logistic regression is not required. Therefore, we need to satisfy the Sukhatme–Randles–deWet condition. Parameter estimates produced by logistic regression and the LDA are both consistent under assumption of normality [35]; therefore, when sample size is sufficiently large, the two estimates are very close. In Table I, we used logistic regression to estimate parameters for



**Figure 2.** Two normal density curves with empirical (dotted line) and theoretical from Table II (solid line) variances overlaid on the histograms of the five statistics calculated for nested models under the alternative (non-degenerate case). Simulated two predictor variables and binary outcomes:  $x_1, x_2 \mid D = 1 \sim N(\mu, \Sigma)$  and  $x_1, x_2 \mid D = 0 \sim N(0, \Sigma)$ .  $x_1, x_2$  are predictors from the full model;  $x_1$  is the predictor from the reduced model.



simulated normal data. Table I supports the theoretical findings of Statement 2.C and 2.D despite the use of logistic regression.

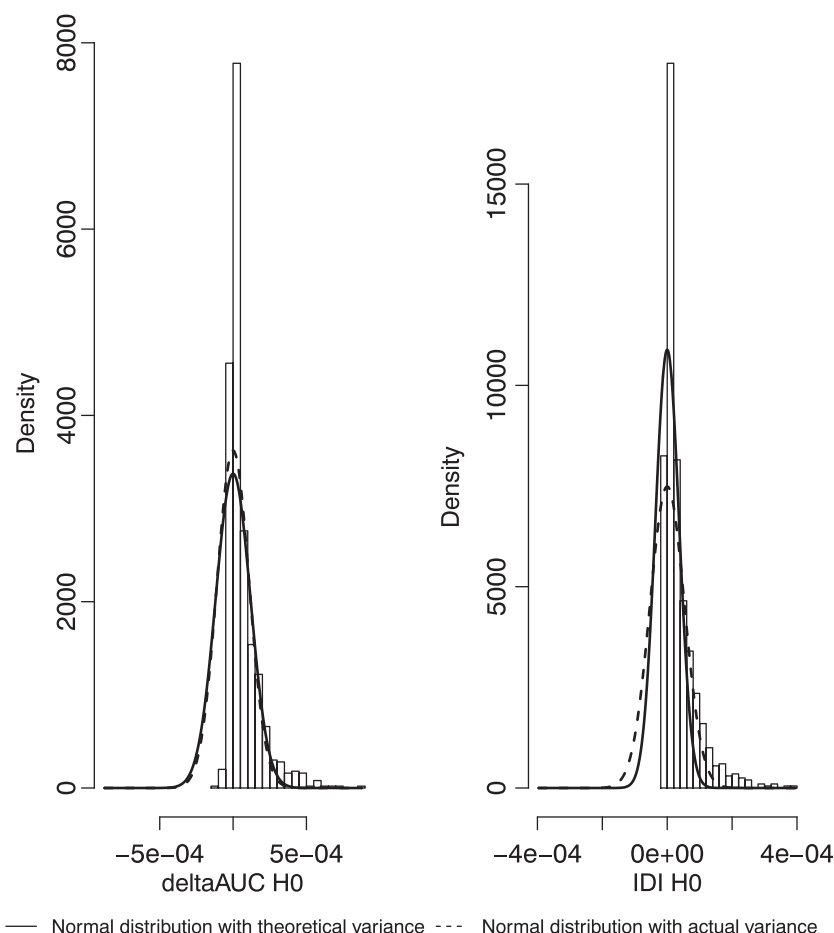
The proof of Statement 2.C and the discussion earlier rely on normality of predictor variables. An important question is how sensitive these results are to the normality assumption. In Section 4, we apply the results of this section to real-life non-normal data using logistic regression and discuss the implications.

### 3.6. Degenerate case

In the Appendix, we show that when comparing nested models under the null,  $\Delta\text{AUC}$ , NRIs, and IDI belong to a degenerate class of U-statistics. They are distributed as an infinite sum of weighted Chi-square distributions. Histograms in Figure 3 demonstrate why any test that assumes normality is invalid for  $\Delta\text{AUC}$  and IDI.

### 3.7. Injecting random noise to remedy degeneracy

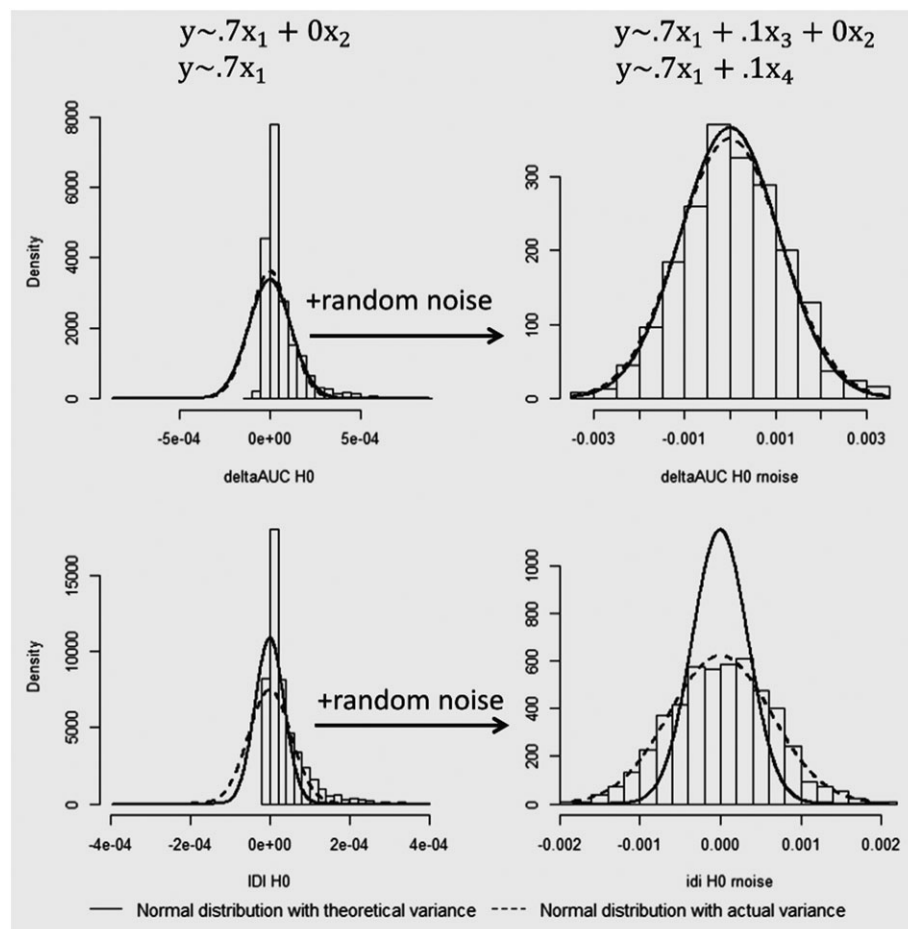
In previous sections, we discussed problems induced by the degenerate state of  $\Delta\text{AUC}$ , NRIs, and IDI when they compare nested models under the null. Their asymptotic distribution and variance estimators become practically intractable. In their non-degenerate state,  $\Delta\text{AUC}$ , NRIs, and IDI follow a normal distribution asymptotically, and variance formulas are available. In this section, we show how degeneracy is at the root of the problem. We will artificially move  $\Delta\text{AUC}$ , NRIs, and IDI away from degeneracy and show that their distribution functions shift to normal distribution. This will shed some



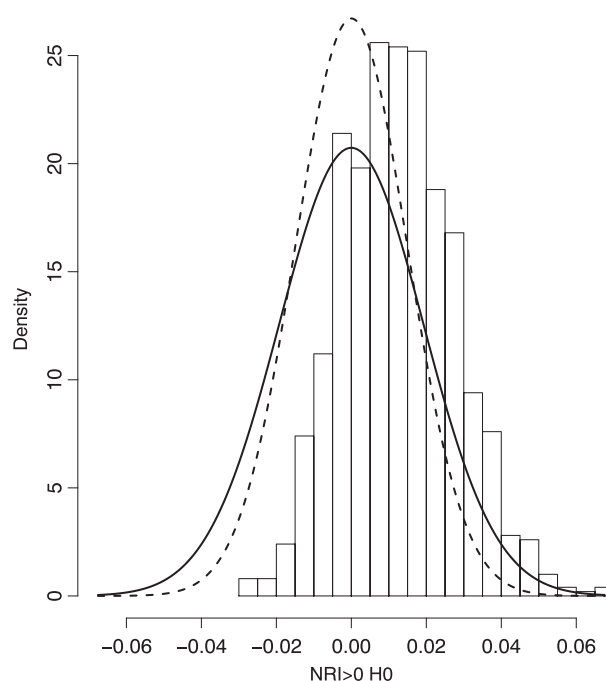
**Figure 3.** Histograms of  $\Delta\text{AUC}$  and IDI when comparing nested models under the null. Two normal density curves with empirical (dotted line) and theoretical (solid line) variances overlaid on the histograms of  $\Delta\text{AUC}$  and IDI calculated for nested models under the null (degenerate case). Simulated two predictor variables and binary outcome:  $x_1 | D = 1 \sim N(\mu, \sigma^2)$ ,  $x_1 | D = 0 \sim N(0, \sigma^2)$ , and  $x_2 \sim N(0, \sigma^2)$ .  $x_2$  is an uninformative predictor.  $x_1, x_2$  are predictors from the full model;  $x_1$  is the predictor from the reduced model.

light on other aspects of NRI behavior that we will discuss later in the section. In the Appendix, the degeneracy condition is formulated in mathematical terms, and it follows that the nested models under the null are the fundamental reason for the degeneracy of  $\Delta\text{AUC}$ ,  $\text{NRI}_{>0}$ , IDI, and all categorical versions of the NRI. So let us consider two nested models under the null.  $\Delta\text{AUC}$ , NRIs, and IDI calculated for these two models will be in a degenerate state. To force them to move away from the degeneracy, we need to violate the degeneracy condition: one way is to force the models away from the null, and an alternative way is to un-nest them. In practical situations, we have no control over a model being under the null or under the alternative. However, we can try to un-nest the two models by injecting random noise, that is, add a weak predictor to the smaller model and another independent weak predictor of the same strength to the other model. Histograms of these statistics for the same models but with injected noise are in the right column of Figure 4. Their distributions shift to asymptotic normality. Results for variance estimators hold in this example too: Variance estimate of  $\Delta\text{AUC}$  is still satisfactory, and the variance of the IDI is underestimated by existing formulas. Our simulations indicate that de-degenerating these two U-statistics comes at the price of a substantial increase of variance and leads to a loss of power. However, this exercise helps to explain why the distribution of  $\text{NRI}_{>0}$ ,  $\text{NRI}(r)$ , and  $3\text{cNRI}$  appears more Gaussian for the degenerate state in our simulations (Figure 5).

The IDI can be written as follows:  $\text{IDI} = \frac{\sum_{i=1}^{n_1} P_{\text{new ev}} - P_{\text{old ev}}}{n_1} - \frac{\sum_{i=1}^{n_0} P_{\text{new nonev}} - P_{\text{old nonev}}}{n_0}$ . The  $\text{NRI}_{>0}$  uses the same definition as IDI but dichotomizes the change in predictive probability:



**Figure 4.** Left column: two nested models under the  $H_0$ ; right column: the two models after un-nesting, preserving the  $H_0$ . Left panel models:  $x_1, x_2$  are predictors from the full model;  $x_1$  is predictor from the reduced model.  $x_1 | D = 1 \sim N(\mu, \sigma^2)$ ,  $x_1 | D = 0 \sim N(0, \sigma^2)$ , and  $x_2 \sim N(0, \sigma^2)$ .  $x_2$  is an uninformative predictor. Right panel models:  $x_1, x_2, x_3$  are predictors from the full model;  $x_1, x_4$  are predictors from the reduced model.  $x_1 | D = 1 \sim N(\mu, \sigma^2)$ ,  $x_1 | D = 0 \sim N(0, \sigma^2)$ , and  $x_2 \sim N(0, \sigma^2)$ .  $x_{3,4} | D = 1 \sim N(\epsilon, I\sigma^2)$  and  $x_{3,4} | D = 0 \sim N(0, I\sigma^2)$ .  $x_2$  is an uninformative predictor, and  $x_3, x_4$  are added 'noise' – independent simulated weak predictors.



**Figure 5.** Histogram of  $NRI > 0$  under  $H_0$ . Simulated two predictor variables and binary outcomes:  $x_1 | D = 1 \sim N(\mu, \sigma^2)$ ,  $x_1 | D = 0 \sim N(0, \sigma^2)$ , and  $x_2 \sim N(0, \sigma^2)$ .  $x_2$  is an uninformative predictor.  $x_1, x_2$  are predictors from the full model;  $x_1$  is the predictor from the reduced model.

$$NRI_{>0} = \frac{\sum_{i=1}^{n_1} \text{Sign}[p_{\text{new ev}} - p_{\text{old ev}}]}{n_1} - \frac{\sum_{i=1}^{n_0} \text{Sign}[p_{\text{new nonev}} - p_{\text{old nonev}}]}{n_0}$$

Therefore, we can view the  $NRI_{>0}$  as an IDI that adds to each summand a random component that complements it to the nearest of the values of 1 or  $-1$ . This random component operates as injected noise in Figure 4. It adds enough noise so that  $NRI_{>0}$  transitions to non-degeneracy and its histogram looks Gaussian, even when adding a predictor variable of interest ( $x_2$ ) does not improve the performance of the model (Figure 5). Note that  $NRI_{>0}$  remains biased. Its bias is studied in [39].

#### 4. Practical example

We apply our results to FHS [1,40] data. Full information about this data set and the study including the enrollment criteria is reported in [40]. Briefly, 8365 people free of cardiovascular disease at baseline examination were followed for 12 years. The outcome of interest was coronary heart disease (CHD), and 640 people developed CHD during follow-up (7.7%). Predictor variables in this example include age, total and high-density lipoprotein cholesterol, systolic and diastolic blood pressure, baseline diabetes status, and current smoking. All continuous variables are log-transformed. We use logistic regression to run the full model with all the predictors. We also ran a series of smaller nested models, which we obtained by omitting from the full model one of the predictor variables.

The bootstrap estimator of the standard error is consistent for a wide range of statistics under mild regularity conditions [41–43]. Therefore, we can use the bootstrap estimate of the standard error of  $\Delta AUC$ ,  $NRI_{>0}$ , and IDI as a proxy for the gold standard, that is, as an estimator with established consistency. For this reason, we define the relative bias of the formula-based standard error as the difference between the average of a formula-based and bootstrap-based variance estimates divided by the bootstrap-based variance estimate.

In this practical example, all predictors are statistically significant; therefore according to results of this paper,  $\Delta AUC$ ,  $NRI_{>0}$ , and IDI are non-degenerate U-statistics, and according to Statement 2.C,

we would expect low bias of the theoretical standard error formulas for  $\Delta\text{AUC}$ ,  $\text{NRI}_{>0}$ , and  $\text{NRI}(\text{r})$  and high bias for those that require adjustment for estimated parameters:  $3\text{cNRI}$  and  $\text{IDI}$ .

## 5. Results

Relative bias of the standard error was calculated for FHS data using bootstrap as described in the previous section. Results are presented in Table IV.

As we anticipate, the two statistics that require adjustment for estimated parameters ( $\text{IDI}$  and  $3\text{cNRI}$ ) have a stable strong bias in Table III. However, contrary to our expectations, bias of the theoretical standard error estimates of the three statistics that should not require adjustment for estimated parameters ( $\Delta\text{AUC}$ ,  $\text{NRI}_{>0}$ , and  $\text{NRI}(\text{r})$ ) varies greatly. For example, the DeLong formula for standard error of  $\Delta\text{AUC}$  often underestimates it by as much as 23% and the formula for  $\text{NRI}_{>0}$  by as much as 56%. Statement 2 is proved under assumption of normally distributed predictors, and this result is consistent with empirical simulations in Table I. But some of the simulations with real-life data in Table IV still show substantial bias. To further explore this phenomenon, we first check the stability of our results in Table III. We replicate bootstrap analysis several times with the FHS data set but with different random seed. Relative bias remains present across replications. Second, we use the result obtained by Harrell *et al.* [21], that is, that tests of c-index (a survival analysis version of AUC) have very low power. We hypothesize that  $\Delta\text{AUC}$ ,  $\text{NRI}_{>0}$ , and  $\text{NRI}(\text{r})$  experience similar loss of power. We observed in our simulations that transition to non-degeneracy is gradual (Appendix Figure A1), so lack of power may be explained by degenerate behavior of the  $\Delta\text{AUC}$ ,  $\text{NRIs}$ , and  $\text{IDI}$  even for moderately strong predictor variables; therefore, we cannot use standard error formulas developed under the assumption of non-degeneracy. This reasoning implies that if we artificially inflate the strength of the added predictor variable,  $\Delta\text{AUC}$ ,  $\text{NRI}_{>0}$ , and  $\text{NRI}(\text{r})$  should move further away from the null and the relative bias of standard error estimates of  $\Delta\text{AUC}$ ,  $\text{NRI}_{>0}$ , and  $\text{NRI}(\text{r})$  will go down. Standard error estimates of  $3\text{cNRI}$  and  $\text{IDI}$  have another problem: they require adjustment for estimated parameters. This problem cannot be solved by artificial inflation of effect size, so we expect bias of their standard error estimates

**Table IV.** Relative bias (bottom table) ( $\text{rel.bias} = \frac{\text{SE}_{\text{formula-based}} - \text{SE}_{\text{bootstrap}}}{\text{SE}_{\text{bootstrap}}} \times 100\%$ ) and averaged bootstrap estimates were calculated by bootstrapping FHS dataset by simple random sampling without replacement. The full model included baseline values of age, TCL, HDL, SBP, DBP, diabetes status, and current smoking. The first row compares the full model to the same model without SBP as a predictor. The row 'AGE' compares the full model to the same model but with age omitted from the list of predictors.

	Parameter estimates					
	$\beta^a$	$\Delta\text{AUC}$	$\text{NRI}_{>0}$	$\text{NRI}(\text{r})$	$3\text{cNRI}^c$	$\text{IDI}$
SBP	1.39	0.00	0.03	0.01	0.00	0.00
HDL	-1.65	0.03	0.44	0.06	0.12	0.02
TCL	1.57	0.01	0.25	0.02	0.06	0.01
AGE	2.73	0.02	0.46	0.06	0.16	0.01
DBP	1.05	0.00	0.20	0.01	0.01	0.00
SMOKING	0.50	0.01	0.16	0.02	0.05	0.00
DIABETES	0.55	0.00	-0.10	0.00	0.00	0.00
	Relative bias (%) of standard error estimate of					
	z-score( $\beta$ ) <sup>b</sup>	$\Delta\text{AUC}$	$\text{NRI}_{>0}$	$\text{NRI}(\text{r})$	$3\text{cNRI}^c$	$\text{IDI}$
SBP	3.05	-23	-3	-24	-24	-45
HDL	-10.93	-2	-6	-8	-28	-38
TOT	6.70	-2	-4	-13	-26	-29
AGE	10.41	9	-5	-6	-34	-36
DBP	2.16	20	-36	-19	-36	-16
SMOKING	5.62	-13	-1	-16	-29	-27
DIABETES	3.61	-3	-56	-16	-19	-43

<sup>a</sup> $\beta$  is the linear coefficient by added predictor variable (larger model). Continuous predictors were log-transformed but not standardized.

<sup>b</sup>z-scores of  $\beta$  coefficients ( $\beta/\text{se}(\beta)$ ).

<sup>c</sup>About 2% and 10% cutoffs were used to calculate  $3\text{cNRI}$ .

to stay strong. Table V shows the results of the bootstrap for the same data as in Table IV, but with artificially inflated effect sizes of added predictor variables.

Because we have artificially forced predictor variables away from the null, results presented in Table V now support Statement 2. As expected, formula-based standard error estimates of  $\Delta\text{AUC}$ ,  $\text{NRI}_{>0}$ , and  $\text{NRI}(\text{r})$  have low bias, and 3cNRI and IDI have high bias because the latter group requires adjustment for estimated parameters.

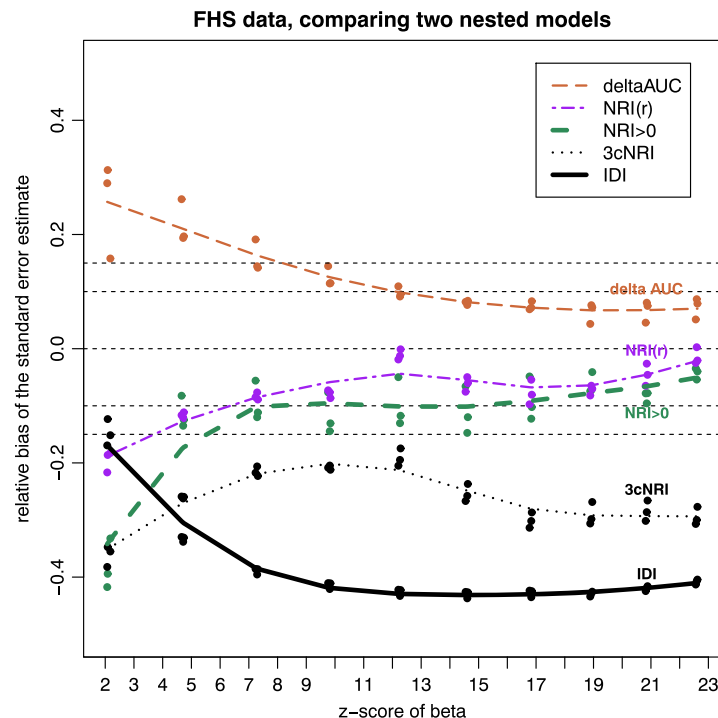
In Figure 6, we illustrate the relationship between relative bias of formula-based standard error and effect size of the added predictor.

Results of this bootstrap analysis using real-life data suggest that Statement 2 is sensitive to the assumption of non-degeneracy. Statistical significance at the 0.05 level of added predictor variable is not sufficient to guarantee non-degeneracy, and associations with stronger effect sizes are required for asymptotic formulas to become consistent. In our example, when  $p$ -values of added predictor variables are weaker than  $10^{-5}$ ,  $\Delta\text{AUC}$ ,  $\text{NRI}_{>0}$ , and  $\text{NRI}(\text{r})$  are too close to degeneracy. Figure A1 in the Appendix illustrates very slow gradual transition away from degeneracy of  $\Delta\text{AUC}$  as the added predictor variable gets stronger. That is, the distribution of  $\Delta\text{AUC}$  is still non-normal when the  $z$ -score of the added predictor variable is less than 4.0 ( $p\text{-value} \leq 6 \cdot 10^{-5}$ ). Much stronger effect sizes are needed to achieve non-degeneracy. This observation explains why formula-based standard error estimators of  $\text{NRI}_{>0}$ ,  $\text{NRI}(\text{r})$ , and  $\Delta\text{AUC}$  are biased in Table IV when  $p$ -values of the added predictor variable are less than 0.05 but greater than  $10^{-5}$ .

Figure 6 illustrates that in FHS data, as the added predictor variable gets stronger, bias of standard error estimates of  $\Delta\text{AUC}$ ,  $\text{NRI}_{>0}$ , and  $\text{NRI}(\text{r})$  decreases. With  $z$ -scores of beta coefficient  $\geq 4.0$  relative bias of formula-based standard error estimates of  $\text{NRI}_{>0}$ ,  $\text{NRI}(\text{r})$  falls below 15% while standard error of  $\Delta\text{AUC}$  is still overestimated by the asymptotic formula and requires an even stronger predictor to lower its relative bias below 15%. When the  $z$ -score of the added predictor in nested models framework is less than 4.0 ( $p\text{-value} > 10^{-5}$ ), standard errors of  $\Delta\text{AUC}$ ,  $\text{NRI}_{>0}$ , and  $\text{NRI}(\text{r})$  should be estimated using resampling methods. Electronic Health Records, pooled genetic cohorts, social networks data, etc. can result in very large sample sizes and potentially very low  $p$ -values. For such large sample sizes, traditional resampling technique can become time consuming. Our results show that in this situation, formula-based standard error estimates of  $\Delta\text{AUC}$ ,  $\text{NRI}_{>0}$ , and  $\text{NRI}(\text{r})$  may have low bias and may be estimated by using the formulas presented in Table II. Table V implies that bias of added dichotomous predictors may remain strong in all scenarios. Standard errors of 3cNRI and IDI always require adjustment for estimated

**Table V.** Analysis presented in Table III was repeated, but effect size of each added predictor variable was artificially inflated.  $\text{rel.bias} = \frac{\text{SE}_{\text{formula-based}} - \text{SE}_{\text{bootstrap}}}{\text{SE}_{\text{bootstrap}}} \times 100\%$ . Effect size of dichotomous variables was inflated by artificially increasing their prevalence among events.

	Parameters Estimates					
	$\beta$	$\Delta\text{AUC}$	$\text{NRI}_{>0}$	$\text{NRI}(\text{r})$	3cNRI	IDI
SBP	21.70	0.20	1.50	0.44	0.95	0.43
HDL	5.63	0.17	1.16	0.34	0.70	0.26
TCL	14.35	0.21	1.63	0.48	1.08	0.55
AGE	2.72	0.02	0.46	0.06	0.16	0.01
DBP	36.90	0.22	1.83	0.55	1.23	0.77
SMOKING	1.34	0.04	0.56	0.09	0.18	0.04
DIABETES	2.32	0.06	0.62	0.06	0.12	0.11
	Relative Bias (%) of Standard Error Estimate of					
	$z\text{score}(\beta)$	$\Delta\text{AUC}$	$\text{NRI}_{>0}$	$\text{NRI}(\text{r})$	3cNRI	IDI
SBP	29.46	-2	-12	0	-36	-26
HDL	27.46	2	-7	-2	-37	-25
TCL	27.90	0	-7	-7	-38	-22
AGE	10.41	7	-4	-8	-35	-35
DBP	22.58	1	-9	-2	-40	-27
SMOKING	14.63	-2	-2	-3	-24	-39
DIABETES	21.34	0.02	-0.23	0.08	-0.11	-0.38



**Figure 6.** Relative bias of standard error estimate as a function of strength of the added predictor variable (z-score of  $\beta_{DBP}$ ) using Framingham Heart Study data. Reduced model included predictor variables: age, HDL and total cholesterol, systolic blood pressure, smoking, and diabetes status. Full model = reduced model + diastolic blood pressure (DBP). We artificially varied the strength of added DBP variable and calculated relative bias of variance estimate using theoretical formula relative to its bootstrapped value.  $z\text{-score}(\beta_{DBP}) = \beta_{DBP}/se(\beta_{DBP})$ .  $rel.bias = (se_{formula,based} - se_{bootstrap})/se_{bootstrap}$ . [Colour figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

parameters. As illustrated in Figure 6, their bias stays strong. For these reasons, resampling methods should be preferred in all situations to estimate standard errors of 3cNRI and IDI. We recommend similar strategies in estimating confidence intervals for  $\Delta AUC$ ,  $NRI_{>0}$ ,  $NRI(r)$ , 3cNRI, and IDI.

## 6. Discussion

This validation study shows that the behavior of  $\Delta AUC$ ,  $NRI_{>0}$ ,  $NRI(r)$ , 3cNRI, and IDI is affected by the interplay of several factors including the shift to degeneracy (non-normality) when comparing two nested models under the null and the lack of adjustment for estimated parameters for 3cNRI and IDI.

Our results explicitly specify conditions under which normal distribution theory can and cannot be applied to  $\Delta AUC$ ,  $NRI_{>0}$ , IDI, and categorical versions of the NRI when comparing two nested models. A few tests of these statistics have been proposed, and all with the exception of [20] rely on asymptotic normality. Our results imply that tests that rely on asymptotic normality are invalid for nested models and should not be used. Fortunately, testing is unnecessary: Pepe *et al.* [26] proved that testing of several measures of model performance is redundant because improvement in most of these statistics is equivalent to the significance of the new predictor variable. Therefore, the recommended strategy is to establish the significance of the regression coefficient first and then evaluate improvement in model performance by producing confidence intervals for measures of performance such as  $\Delta AUC$ , NRIs, and IDI.

Using U-statistics theory, we proved that when the added predictor variable is significant, the distribution of  $\Delta AUC$ ,  $NRI_{>0}$ ,  $NRI(r)$ , 3cNRI, and IDI is normal; therefore, asymptotic confidence intervals can rely on the normal distribution. We considered their variance estimators and showed in Statement 2 that theoretical standard error estimates of  $\Delta AUC$ ,  $NRI_{>0}$ , and  $NRI(r)$  are valid when predictor variables are normally distributed. Our practical example using FHS data demonstrated that when the added predictor is significant but the  $p$ -value is not particularly low, the variance of  $NRI_{>0}$



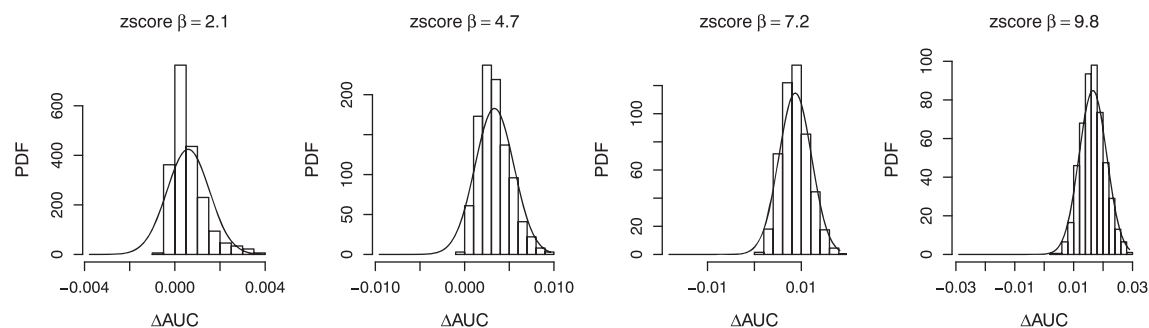
and  $\text{NRI}(r)$  is still underestimated by the formula and the variance of  $\Delta\text{AUC}$  is overestimated. Our simulations demonstrated that a stronger added predictor variable is required to reach non-degeneracy, a necessary condition for validity of the formulas. We offer an example in which the  $p$ -value of added predictor variable  $<10^{-5}$  is needed for  $\Delta\text{AUC}$ ,  $\text{NRI}_{>0}$ , and  $\text{NRI}(r)$  to fully transition to non-degeneracy and for the relative bias of the standard error of  $\text{NRI}_{>0}$  and  $\text{NRI}(r)$  to drop to below 15% (Figure 6). Such high effect sizes and significance levels might be common in Big Data studies.

While formula-based standard errors of  $\Delta\text{AUC}$ ,  $\text{NRI}_{>0}$ , and  $\text{NRI}(r)$  are valid in the situations described previously, formula-based standard error estimators of 3cNRI and IDI are not. Unless they are adjusted for estimated parameters, they underestimate actual variance. Therefore, existing standard errors formulas for 3cNRI and IDI should not be used, and bootstrap or other resampling technique should be employed instead.

Additionally, using U-statistics theory, we showed that the standard error estimator of  $\Delta\text{AUC}$  can be calculated as the variance of the change in ranks of predicted probabilities (Table III). In our numerical simulations, the new variance estimator was identical to the one produced by DeLong *et al.* [11] and the two are likely algebraically equivalent when there are no ties in predicted probabilities. However, rigorous proof of this result is beyond the scope of this paper.

In summary, when comparing two nested models after establishing the significance of the regression coefficient of an added predictor variable, we recommend estimating formula-based standard errors and confidence intervals of  $\Delta\text{AUC}$ ,  $\text{NRI}_{>0}$ , and  $\text{NRI}(r)$  when the significance of predictor variables is strong enough ( $p$ -value  $<10^{-5}$ ,  $z$ -score  $>4.0$  in our FHS data example). In other situations, the CIs of  $\Delta\text{AUC}$  are too conservative; while CIs for all other statistics are too narrow, therefore, resampling techniques (such as bootstrap) should be used to estimate these. Standard errors of IDI and 3-category NRI should always be estimated by the bootstrap or other resampling technique.

## Appendix A



**Figure A1.** Transition of  $\Delta\text{AUC}$  from degeneracy to non-degeneracy. Framingham Heart Study data. Reduced model included predictor variables: age, high-density lipoprotein and total cholesterol, systolic blood pressure, smoking, and diabetes status. Full model = reduced model + diastolic blood pressure (DPF). We artificially varied the strength of added DPF variable. Strength is measured by  $z\text{-score}(\beta_{\text{dpf}}) = \beta_{\text{dpf}}/\text{se}(\beta_{\text{dpf}})$ .

### A.1. Brief review of U-statistics theory definitions and results

U-statistics theory can be viewed as an extension of Central Limit Theorem for sums of correlated variables [1]. Generalized U-statistics can be defined for a sample generated from  $k$  different distributions with no restriction on  $k$ , but for our purposes, we focus on the two-sample generalized U-statistics.

#### Definition 1

*Generalized two-sample U-statistics* [1]. Define  $x_1, \dots, x_{n_0}$  and  $y_1, \dots, y_{n_1}$  as samples of size  $n_0$  and  $n_1$ , respectively, from two different (possibly multivariate) distribution functions  $\mathbf{x} \sim \mathbf{F}(\cdot)$  and  $\mathbf{y} \sim \mathbf{G}(\cdot)$ .

Assume that  $x_1, \dots, x_{n_0}$  and  $y_1, \dots, y_{n_1}$  are independent within their samples and also across the two samples. Then *generalized U-statistic* is defined as follows:

$$U_{n_0, n_1} = \binom{n_0}{c}^{-1} \binom{n_1}{d}^{-1} \sum_{(n_0, c)} \sum_{(n_1, d)} \psi(x_{i_1}, \dots, x_{i_c}; y_{j_1}, \dots, y_{j_d}), \quad (a1)$$

where  $\psi(\cdot, \cdot)$  is a real-valued kernel symmetric in each set of arguments and  $c$  and  $d$  are dimension constants satisfying  $1 \leq c \leq n_0, 1 \leq d \leq n_1$ . Note that  $c = d = 1$  for  $\Delta\text{AUC}$ ,  $\text{NRI} > 0$ , and  $\text{IDI}$  (will be demonstrated in the following). Therefore, it is enough to consider  $c = d = 1$ . In this case, the symmetry condition is automatically satisfied and formula (a1) simplifies to

$$U_{n_0, n_1} = \frac{1}{n_0 n_1} \sum_{i=1}^{n_0} \sum_{j=1}^{n_1} \psi(x_i; y_j) \quad (a2)$$

U-statistics can be viewed as a sum of non-i.i.d. (where i.i.d. is independent and identically distributed) random variables for which Central Limit Theorem can be extended

*Theorem [1]*

$\sqrt{N} (U_{n_0, n_1} - \Theta) \xrightarrow{D} \mathcal{N}(0, \delta^2), N = n_0 + n_1 \rightarrow \infty$ , unless  $U_{n_0, n_1}$  is a degenerate U-statistic.

See Lee [1] for a definition of  $\delta^2$ .

Therefore, this Theorem fully defines asymptotic distribution of non-degenerate U-statistics. If we show that  $\Delta\text{AUC}$ ,  $\text{NRI} > 0$ , and  $\text{IDI}$  are non-degenerate U-statistics, we can use this Theorem in order to create their tests and confidence intervals. Degeneracy condition is defined in the following.

#### A.2. Brief review of degeneracy definitions for U-statistics

To utilize the Theorem, we need first to check degeneracy condition and then calculate variance; both of them are defined in terms of H-decomposition of U-statistics. Therefore, we first introduce H-decomposition and formulate degeneracy condition for any U-statistics and then calculate H-decomposition and define degeneracy condition of  $\Delta\text{AUC}$ ,  $\text{NRI}_{>0}$ , and  $\text{IDI}$  in the following.

##### Definition 2

First recall that  $x_1, \dots, x_{n_0}$  and  $y_1, \dots, y_{n_1}$  are samples of size  $n_0$  and  $n_1$ , respectively, from two different (possibly multivariate) distribution functions  $\mathbf{x} \sim \mathbf{F}(\cdot)$  and  $\mathbf{y} \sim \mathbf{G}(\cdot)$ .

*H-decomposition of a generalized two-sample U-statistics [1–2].* Any generalized two-sample U-statistic (defined in (a2)) where  $\mathbf{x} \sim \mathbf{F}(\cdot)$  and  $\mathbf{y} \sim \mathbf{G}(\cdot)$  can be decomposed into a summation of several sums of i.i.d. random variables as follows:

$$\begin{aligned} U_{n_0, n_1} &\stackrel{\text{def}}{=} \frac{1}{n_0 n_1} \sum_{i=1}^{n_0} \sum_{j=1}^{n_1} \psi(x_i; y_j) \\ &= h^{(0, 0)} + \frac{1}{n_1} \sum_{i=1}^{n_0} h^{(1, 0)}(x_i) + \frac{1}{n_1} \sum_{j=1}^{n_1} h^{(0, 1)}(y_j) + \frac{1}{n_0 n_1} \sum_{i=1}^{n_0} \sum_{j=1}^{n_1} h^{(1, 1)}(x_i; y_j), \end{aligned} \quad (a3)$$

where  $h^{(\cdot, \cdot)}(\cdot)$  are real-valued kernels. Formulas to calculate the kernels  $h^{(\cdot, \cdot)}(\cdot)$  are given by Lee [1]:

$$\begin{aligned} h^{(0, 0)} &= \iint \psi(x; y) d\mathbf{F}(x) d\mathbf{G}(y) \\ h^{(1, 0)}(x_i) &= \int \psi(x_i; y) d\mathbf{G}(y) - \iint \psi(x; y) d\mathbf{F}(x) d\mathbf{G}(y) \\ h^{(0, 1)}(y_j) &= \int \psi(x; y_j) d\mathbf{F}(x) - \iint \psi(x; y) d\mathbf{F}(x) d\mathbf{G}(y) \\ h^{(1, 1)}(x_i; y_j) &= \psi(x_i; y_j) - \int \psi(x_i; y) d\mathbf{G}(y) - \int \psi(x; y_j) d\mathbf{F}(x) + \iint \psi(x; y) d\mathbf{F}(x) d\mathbf{G}(y) \end{aligned} \quad (a4)$$

where  $\mathbf{F}(\cdot)$  and  $\mathbf{G}(\cdot)$  are cumulative distribution functions of random vectors  $\mathbf{x}$  and  $\mathbf{y}$ ,  $\mathbf{x} \sim \mathbf{F}(\cdot)$  and  $\mathbf{y} \sim \mathbf{G}(\cdot)$ .

##### Definition 3

*Degeneracy of generalized U-statistics.* Using the previous notation, if  $\text{Var}(h^{(0, 1)}(\mathbf{y})) = 0$  and  $\text{Var}(h^{(1, 0)}(\mathbf{x})) = 0$  for some family of distribution functions  $\mathbf{F}(\cdot) \in \Phi(\cdot)$  and  $\mathbf{G}(\cdot) \in \Gamma(\cdot)$ , then generalized U-statistic H-decomposed as in Definition 2 is called degenerate.

It follows that in order to check degeneracy, it is enough to calculate the first two terms  $h^{(1,0)}(x_i)$ ,  $h^{(0,1)}(y_j)$  of H-decomposition. Let us calculate  $h^{(1,0)}(x_i)$ ,  $h^{(0,1)}(y_j)$  for  $\Delta AUC$ ,  $NRI > 0$ ,  $IDI$ , and  $3cNRI$  and check their degeneracy conditions.

For two nested models,  $\Delta AUC$ ,  $NRI_{>0}$ ,  $3cNRI$ , and  $IDI$

1. are generalized U-statistics.
2. are degenerate generalized U-statistics if models are under the null. If they are degenerate, they do not follow normal distribution.
3. are non-degenerate generalized U-statistics under the alternative and therefore are asymptotically normally distributed.
4. If we do not adjust for estimated parameters, then the variance formulas derived from U-statistics theory are same as the published formulas [3–4].

For two non-nested models,  $\Delta AUC$ ,  $NRI > 0$ ,  $3cNRI$ , and  $IDI$

1. are generalized U-statistics.
2. are non-degenerate generalized U-statistics.
3. are asymptotically normally distributed with the variance formulas provided in the following.

Variance formulas:

$$\begin{aligned}\hat{\sigma}_{\Delta AUC}^2 &= \frac{\text{Var}(\text{rank}_e^*(a^T x_i) - \text{rank}_e(a^T x_i))}{n_0} + \frac{\text{Var}(\text{rank}_{ne}^*(a^T y_j) - \text{rank}_{ne}(a^T y_j))}{n_1}, \\ \hat{\sigma}_{NRI}^2 &= \frac{\hat{p}_{ne}^{up}(1 - \hat{p}_{ne}^{up})}{n_0} + \frac{\hat{p}_e^{up}(1 - \hat{p}_e^{up})}{n_1}, \\ \hat{\sigma}_{3cNRI}^2 &= \frac{4(\hat{p}_{ne}^{2up} + \hat{p}_{ne}^{2down}) + \hat{p}_{ne}^{1up} + \hat{p}_{ne}^{1down} - (2(\hat{p}_{ne}^{2up} - \hat{p}_{ne}^{2down}) + \hat{p}_{ne}^{1up} - \hat{p}_{ne}^{1down})^2}{n_0} \\ &\quad + \frac{4(\hat{p}_{ev}^{2up} + \hat{p}_{ev}^{2down}) + \hat{p}_{ev}^{1up} + \hat{p}_{ev}^{1down} - (2(\hat{p}_{ev}^{2up} - \hat{p}_{ev}^{2down}) + \hat{p}_{ev}^{1up} - \hat{p}_{ev}^{1down})^2}{n_1}, \\ \hat{\sigma}_{IDI}^2 &= \frac{\text{Var}(\Delta \text{predp}(x_i))}{n_0} + \frac{\text{Var}(\Delta \text{predp}(y_j))}{n_1},\end{aligned}$$

where

$\text{rank}_e(a^T x_i)$  is the rank of  $a^T x_i$  among  $a^T y$  (risk scores of events),  
 $\text{rank}_{ne}(a^T y_j)$  is the rank of  $a^T y_j$  among  $a^T x_i$  (risk scores of events) and as before \* means that similar functions and coefficients were calculated for the old model.

$$\begin{aligned}\hat{p}_{ne}^{up} &= \frac{\#non - events \text{ moving up}}{n_0}, \\ \hat{p}_e^{up} &= \frac{\#events \text{ moving up}}{n_1},\end{aligned}$$

$\Delta \text{predp}(x_i)$  is the difference in predicted probabilities from two models for non-event  $x_i$

$\Delta \text{predp}(y_j)$  is the difference in predicted probabilities from two models for event  $y_j$ .

## Appendix B: Proof of the main result

### B.1. Main result

$\Delta AUC$ ,  $NRI_{>0}$ ,  $NRI(r)$ ,  $3cNRI$  and  $IDI$

STATEMENT 1 are **generalized U-statistics with estimated parameters**.

STATEMENT 2 belong to **non-degenerate subclass** if they compare any non-nested models or nested models under the alternative. As non-degenerate U-statistics, they

a follow normal distribution asymptotically.

- b Available variance formulas are algebraically equal to the variance estimates provided by U-statistics theory if we ignore adjustment for estimated parameters.
- c Variances of  $\Delta AUC$ ,  $NRI_{>0}$ , and  $NRI(r)$  do not need to be adjusted for estimated parameters if predictor variables are normally distributed.
- d Variance of IDI and 3-category NRI always should be adjusted for estimated parameters.

STATEMENT 3  $\Delta AUC$ ,  $NRI_{>0}$ ,  $NRI(r)$ ,  $3cNRI$ , and IDI belong to **degenerate subclass** if they compare nested models under the null. As degenerate U-statistics, they do not follow normal distribution and available variance estimates do not apply for them.

B.1.1.1. STATEMENT 1.  $\Delta AUC$ ,  $NRI > 0$ , and IDI are generalized two-sample U-statistics.

Lemma 1

$\Delta AUC$ ,  $NRI > 0$ , IDI, and categorical NRIs are generalized two-sample U-statistics.

Proof

We need to show that each of the three statistics can be written in the form (a2).

B.1.1.1.  $\Delta AUC$  is a generalized two-sample U-statistics. We assume that AUC is estimated by the Mann–Whitney statistic [5] – a non-parametric unbiased estimator, often referred to as the *c*-statistic [6–7]. Mann–Whitney statistic for the full model is as follows:

$$AUC = \frac{1}{n_0 n_1} \sum_{i=1}^{n_0} \sum_{j=1}^{n_1} I[a' x_i < a' y_j],$$

where  $I[\cdot]$  is the indicator function.

AUC for the reduced model is as follows:

$$AUC^* = \frac{1}{n_0 n_1} \sum_{i=1}^{n_0} \sum_{j=1}^{n_1} I[a^* x_i < a^* y_j]$$

If we write  $\Delta AUC$  as (a2), then we will show that  $\Delta AUC$  is a U-statistic.

$$\begin{aligned} \Delta AUC &= \frac{1}{n_0 n_1} \sum_{i=1}^{n_0} \sum_{j=1}^{n_1} I[a' x_i < a' y_j] - \frac{1}{n_0 n_1} \sum_{i=1}^{n_0} \sum_{j=1}^{n_1} I[a^* x_i < a^* y_j] \\ &= \frac{1}{n_0 n_1} \sum_{i=1}^{n_0} \sum_{j=1}^{n_1} \left( I[a' x_i < a' y_j] - I[a^* x_i < a^* y_j] \right) \end{aligned} \quad (a5)$$

Denoting  $\psi_{\Delta AUC}(x, y) = I[a' x_i < a' y_j] - I[a^* x_i < a^* y_j]$  as the kernel, we see that  $\Delta AUC$  satisfies the properties required in Definition 1 for generalized two-sample U-statistics.

B.1.1.2.  $NRI_{>0}$  is a generalized two-sample U-statistics.  $NRI_{>0}$  is defined as follows:

$$\begin{aligned} NRI_{>0} &= \frac{\#up, events}{n_1} - \frac{\#up, nonevents}{n_0} = \frac{1}{n_1 n_0} (n_0 \cdot \#up, events - n_1 \cdot \#up, nonevents) \\ &= \frac{1}{n_0 n_1} \sum_{i=1}^{n_0} \sum_{j=1}^{n_1} \left( I[(a' - a^*) y_j > 0] - I[(a' - a^*) x_i > 0] \right) \end{aligned} \quad (a6)$$

Denoting kernel as  $\psi_{NRI_{>0}}(x, y) = I[(a' - a^*) y_j > 0] - I[(a' - a^*) x_i > 0]$ , we see that continuous NRI ( $NRI > 0$ ) satisfies the properties required in Definition 1 for generalized two-sample U-statistics.

B.1.1.3. Three-category NRI is a generalized two-sample U-statistics. We can define categories as low (L), intermediate (I), and high (H), defined by cutoffs  $c_1$  and  $c_2$ . Then jump size adjusted NRI is defined as follows:

$$3cNRI = \frac{1}{n_0 n_1} \sum_{i=1}^{n_0} \sum_{j=1}^{n_1} \left( I_{3c}[y_j] - I_{3c}[x_i] \right), \quad (a9)$$

where

$$I_{3c}[x_i] = \begin{cases} 2, & \text{if predicted prob of non - event } i \text{ moved two categories up} \\ 1, & \text{if predicted prob of non - event } i \text{ moved one category up} \\ 0, & \text{if predicted prob of non - event } i \text{ did not change} \\ -1, & \text{if predicted prob of non - event } i \text{ moved one category down} \\ -2, & \text{if predicted prob of non - event } i \text{ moved two categories down} \end{cases}$$

Similar definition for  $I_{3c}[y_j]$ .

Denoting kernel as  $\psi_{3cNRI}(x, y) = I_{3c}[y_j] - I_{3c}[x_i]$ , we see that (a9) satisfies the properties required in Definition 1 for generalized two-sample U-statistics.

*B.1.1.4. IDI is a generalized two-sample U-statistics.* Denoting  $f(\cdot) = \text{inv.logit}(\cdot)$ , IDI is defined as follows:

$$\begin{aligned} IDI &= \frac{\sum_1^{n_1} p_j - p_j^*}{n_1} - \frac{\sum_1^{n_0} p_i - p_i^*}{n_0} = \frac{\sum_1^{n_1} f(a' y_j) - f(a^* y_j)}{n_1} - \frac{\sum_1^{n_0} f(a' x_i) - f(a^* x_i)}{n_0} \\ &= \frac{1}{n_0 n_1} \sum_{i=1}^{n_0} \sum_{j=1}^{n_1} \left( f(a' y_j) - f(a^* y_j) - (f(a' x_i) - f(a^* x_i)) \right) \end{aligned} \quad (a7)$$

Denoting kernel as  $\psi_{IDI}(x, y) = f(a' y_j) - f(a^* y_j) - (f(a' x_i) - f(a^* x_i))$ , we see that IDI satisfies the properties required in Definition 1 for generalized two-sample U-statistics.

Q.E.D.

## B.2. Degeneracy condition of $\Delta AUC$ , NRIs, and IDI

*B.2.1. STATEMENT 2.  $\Delta AUC$ , NRIs, and IDI belong to non-degenerate subclass if they compare any non-nested models or nested models under the alternative.*

*Proof*

First, let us write H-decomposition of  $\Delta AUC$ ,  $NRI > 0$ , IDI, and categorical NRIs.

### B.2.1. H-decomposition and degeneracy condition of $\Delta AUC$ .

$$\begin{aligned} h^{(1,0)}(x_i) &= \int \psi(x_i; y) dG(y) - \iint \psi(x; y) dF(x) dG(y) \\ h^{(0,1)}(y_j) &= \int \psi(x; y_j) dF(x) - \iint \psi(x; y) dF(x) dG(y) \end{aligned} \quad (a9)$$

Recall that kernel of  $\Delta AUC$  is defined as follows:

$$\psi_{\Delta AUC}(x, y) = I[a' x_i < a' y_j] - I[a^* x_i < a^* y_j]$$

Plugging it into (a5), we obtain the following:

$$\begin{aligned} h^{(1,0)}(x_i) &= \int I[a' x_i < a' y] - I[a^* x_i < a^* y] dG(y) - h^{(0,0)} \\ &= \Pr(a^T x_i < a^T y) - \Pr(a^{*T} x_i < a^{*T} y) - h^{(0,0)} \\ &= G^*(a^T x_i) - G(a^T x_i) - h^{(0,0)} \end{aligned} \quad (a10)$$

$$\begin{aligned} h^{(0, 1)}(y_i) &= \int I[a'x < a'y_j] - I[a'^*x < a'^*y_j] dF(x) - h^{(0, 0)} \\ &= \Pr(a^Tx < a^Ty_j) - \Pr(a^*Tx < a^*Ty_j) - h^{(0, 0)} \\ &= F(a^Ty_j) - F^*(a^*Ty_j) - h^{(0, 0)} \end{aligned} \quad (a11)$$

where  $h^{(0, 0)} = \iint \psi_{\Delta AUC}(x, y) dF(x) dG(y)$  in (a5), and it is a constant;  $a^Tx_i \sim F(\cdot)$ ,  $a^*Tx_i \sim F^*(\cdot)$ ,  $a^Ty_i \sim G(\cdot)$ ,  $a^*Ty_i \sim G^*(\cdot)$ .

**B.2.1.1.1. Degeneracy condition of  $\Delta AUC$ .**  $\Delta AUC$  is a degenerate U-statistics if  $\text{Var}(h^{(1, 0)}(x_i)) = 0$  and  $\text{Var}(h^{(0, 1)}(y_j)) = 0$ . From (a6), it means that  $G^*(a^*Tx_i) - G(a^Tx_i) \equiv \text{constant} \forall x_i$  and  $F^*(a^*Ty_j) - F(a^Ty_j) \equiv \text{constant} \forall y_j$ , which can happen only if  $a^* = a$ , which is true for nested model under the null.

**B.2.1.1.2. H-decomposition and degeneracy condition of  $NRI_{>0}$ .** Recall that kernel of  $NRI_{>0}$  is as follows:

$$\psi_{NRI_{>0}}(x, y) = I[(a' - a'^*)y_j > 0] - I[(a' - a'^*)x_i > 0]$$

Plugging it into (a10) we obtain the following:

$$\begin{aligned} h^{(1, 0)}(x_i) &= \int I[(a' - a'^*)y_j > 0] - I[(a' - a'^*)x_i > 0] dG(y) - h^{(0, 0)} \\ &= \Pr((a^T - a^*T)y > 0) - I[(a^T - a^*T)x_i > 0] - h^{(0, 0)} \\ h^{(0, 1)}(y_i) &= \int I[(a' - a'^*)y_j > 0] - I[(a' - a'^*)x_i > 0] dF(x) - h^{(0, 0)} \\ &= I[(a^T - a^*T)y_j > 0] - \Pr((a^T - a^*T)x > 0) - h^{(0, 0)} \end{aligned} \quad (a12)$$

Where  $h^{(0, 0)} = \iint \psi_{NRI}(x, y) dF(x) dG(y)$  in (a5), and it is a constant.

**B.2.1.2.1. Degeneracy condition of  $NRI_{>0}$ .**  $NRI > 0$  is a degenerate U-statistics if  $\text{Var}(h^{(1, 0)}(x_i)) = 0$  and  $\text{Var}(h^{(0, 1)}(y_j)) = 0$ . From (a7), it is equivalent to  $I[(a^T - a^*T)x_i > 0] \equiv \text{constant} \forall x_i$  and  $I[(a^T - a^*T)y_j > 0] \equiv \text{constant} \forall y_j$ , which can happen only if  $a^* = a$ , which is true for nested model under the null.

**B.2.1.3. H-decomposition and degeneracy condition of 3-category NRI.** From (a9),

$$\begin{aligned} \psi_{3cNRI}(x, y) &= I_{3c}[y_j] - I_{3c}[x_i], \\ h^{(1, 0)}(x_i) &= -I_{3c}[x_i] + \text{const} \\ h^{(0, 1)}(y_j) &= I_{3c}[y_j] + \text{const} \end{aligned} \quad (a14)$$

**B.2.1.3.1. Degeneracy condition of three-category NRI.** Three-category NRI is a degenerate U-statistics if  $\text{Var}(h^{(1, 0)}(x_i)) = 0$  and  $\text{Var}(h^{(0, 1)}(y_j)) = 0$ . From (a9), it is equivalent to  $I_{3c}[y_j] \equiv \text{constant} \forall y_j$  and  $I_{3c}[x_i] \equiv \text{constant} \forall x_i$ , which can happen only if  $a^* = a$ , which is true for any nested models under the null.

**B.2.1.4. H-decomposition and degeneracy condition of IDI.** Recall that kernel of IDI is as follows:

$$\psi_{IDI}(x, y) = f(a'y_j) - f(a'^*y_j) - (f(a'x_i) - f(a'^*x_i)), \text{ where } f(\cdot) \text{ is the predicted probability.}$$

Plugging it into (a10), we obtain



$$\begin{aligned}
 h^{(1,0)}(x_i) &= \int f(a'y) - f(a^*y) - (f(a'x_i) - f(a^*x_i))dG(y) - h^{(0,0)} \\
 &= f(a'x_i) - f(a^*x_i) + \text{const} \\
 h^{(1,0)}(x_i) &= \int f(a'y_j) - f(a^*y_j) - (f(a'x) - f(a^*x))dF(x) - h^{(0,0)} \\
 &= f(a'y_j) - f(a^*y_j) + \text{const}
 \end{aligned} \tag{a13}$$

*B.2.1.4.1. Degeneracy condition of IDI.* IDI is a degenerate U-statistics if  $\text{Var}(h^{(1,0)}(x_i))=0$  and  $\text{Var}(h^{(0,1)}(y_j))=0$ . From (a8), it is equivalent to  $I[(a^T - a^{*T})x_i > 0] \equiv \text{constant} \forall x_i$  and  $I[(a^T - a^{*T})y_j > 0] \equiv \text{constant} \forall y_j$ , which can happen only if  $a^* = a$ , which is true for any nested models under the null.

Q.E.D.

So all statistics considered in this paper are degenerate if and only if we compare nested models under the null. When the three statistics are non-degenerate, we can apply the Theorem and calculate their variance.

*B.2.2. STATEMENT 2a.* When  $\Delta\text{AUC}$ ,  $\text{NRIs}$ , and  $\text{IDI}$  belong to non-degenerate subclass, they follow normal distribution asymptotically..

$$\begin{aligned}
 \sqrt{N}(\Delta\text{AUC} - \mathbf{E}[\Delta\text{AUC}]) &\xrightarrow{D} N(0, \mathbf{N}\sigma_{\Delta\text{AUC}}^2), \\
 \sqrt{N}(\text{NRI}_{>0} - \mathbf{E}[\text{NRI}_{>0}]) &\xrightarrow{D} N(0, \mathbf{N}\sigma_{\text{NRI}_{>0}}^2), \\
 \sqrt{N}(\text{IDI} - \mathbf{E}[\text{IDI}]) &\xrightarrow{D} N(0, \mathbf{N}\sigma_{\text{IDI}}^2) \\
 \sqrt{N}(3\text{cNRI} - \mathbf{E}[3\text{cNRI}]) &\xrightarrow{D} N(0, \mathbf{N}\sigma_{3\text{cNRI}}^2), \\
 \sqrt{N}(\text{NRI}(\mathbf{p}) - \mathbf{E}[\text{NRI}(\mathbf{p})]) &\xrightarrow{D} N(0, \mathbf{N}\sigma_{\text{NRI}(\mathbf{p})}^2)
 \end{aligned}$$

To derive  $\sigma_{\Delta\text{AUC}}^2$ ,  $\sigma_{\text{NRI}_{>0}}^2$ ,  $\sigma_{\text{IDI}}^2$ ,  $\sigma_{3\text{cNRI}}^2$  and  $\sigma_{\text{NRI}(\mathbf{p})}^2$ , we will rewrite the Theorem in a more detailed way:

**Theorem** [Lee [1] pg 140].

$\sqrt{N}(U_{n_0, n_1} - \Theta) \xrightarrow{D} \mathcal{N}(0, \delta^2)$  as  $N = n_0 + n_1 \rightarrow \infty$ , unless  $U_{n_0, n_1}$  is a degenerate U-statistic. Above  $\delta^2 = \frac{1}{1-p}\delta_{1,0}^2 + \frac{1}{p}\delta_{0,1}^2$ ,  $p = \lim_{N \rightarrow \infty} \frac{n_1}{N}$

and

$$\begin{aligned}
 \delta_{1,0}^2 &= \text{Var}(h^{(1,0)}(x_i)), \\
 \delta_{0,1}^2 &= \text{Var}(h^{(0,1)}(y_j)).
 \end{aligned}$$

Note that  $\delta^2 = \frac{1}{1-p}\delta_{1,0}^2 + \frac{1}{p}\delta_{0,1}^2 \approx \frac{N}{n_0}\delta_{1,0}^2 + \frac{N}{n_1}\delta_{0,1}^2$

*B.2.2.1. Derivation of  $\sigma_{\Delta\text{AUC}}^2$ .*

$$\hat{\sigma}_{\Delta\text{AUC}}^2 = \frac{\text{Var}(\text{rank}_e^*(a^T x_i) - \text{rank}_e(a^T x_i))}{n_0} + \frac{\text{Var}(\text{rank}_{ne}^*(a^T y_j) - \text{rank}_{ne}(a^T y_j))}{n_1}$$

$$\sigma_{\Delta AUC}^2 = \frac{1}{N} \left( \frac{1}{1-p} \delta_{1,0}^2 + \frac{1}{p} \delta_{0,1}^2 \right) \approx \frac{1}{N} \left( \frac{N}{n_0} \delta_{1,0}^2 + \frac{N}{n_1} \delta_{0,1}^2 \right) = \frac{1}{n_0} \delta_{1,0}^2 + \frac{1}{n_1} \delta_{0,1}^2$$

From (a12),

$$\delta_{1,0}^2 = \text{Var} \left( h^{(1, 0)}(\mathbf{x}_i) \right) = \text{Var} \left( G^*(\mathbf{a}^* \mathbf{T} \mathbf{x}_i) - G(\mathbf{a}^T \mathbf{x}_i) \right)$$

$$\delta_{0,1}^2 = \text{Var} \left( h^{(0, 0)}(\mathbf{y}_j) \right) = \text{Var} \left( G^*(\mathbf{a}^* \mathbf{T} \mathbf{x}_i) - G(\mathbf{a}^T \mathbf{x}_i) \right)$$

$G(\mathbf{a}^T \mathbf{x}_i)$  can be estimated as the rank of  $\mathbf{a}^T \mathbf{x}_i$  among  $\mathbf{a}^T \mathbf{y}$  (risk scores of events).

$$G(\mathbf{a}^T \mathbf{x}_i) \approx \text{rank}_e(\mathbf{a}^T \mathbf{x}_i), \quad G^*(\mathbf{a}^* \mathbf{T} \mathbf{x}_i) \approx \text{rank}_e^*(\mathbf{a}^* \mathbf{T} \mathbf{x}_i)$$

$$F(\mathbf{a}^T \mathbf{y}_j) \approx \text{rank}_{ne}(\mathbf{a}^T \mathbf{y}_j), \quad F^*(\mathbf{a}^* \mathbf{T} \mathbf{y}_j) \approx \text{rank}_{ne}^*(\mathbf{a}^* \mathbf{T} \mathbf{y}_j)$$

Therefore,

$$\delta_{1,0}^2 = \text{Var} \left( h^{(1, 0)}(\mathbf{x}_i) \right) \approx \text{Var} \left( \text{rank}_e^*(\mathbf{a}^* \mathbf{T} \mathbf{x}_i) - \text{rank}_e(\mathbf{a}^T \mathbf{x}_i) \right)$$

$$\delta_{0,1}^2 = \text{Var} \left( h^{(0, 1)}(\mathbf{y}_j) \right) \approx \text{Var} \left( \text{rank}_{ne}^*(\mathbf{a}^* \mathbf{T} \mathbf{y}_j) - \text{rank}_{ne}(\mathbf{a}^T \mathbf{y}_j) \right)$$

$$\hat{\sigma}_{\Delta AUC}^2 = \frac{\text{Var} \left( \text{rank}_e^*(\mathbf{a}^* \mathbf{T} \mathbf{x}_i) - \text{rank}_e(\mathbf{a}^T \mathbf{x}_i) \right)}{n_0} + \frac{\text{Var} \left( \text{rank}_{ne}^*(\mathbf{a}^* \mathbf{T} \mathbf{y}_j) - \text{rank}_{ne}(\mathbf{a}^T \mathbf{y}_j) \right)}{n_1}$$

B.2.2.2. Derivation of  $\sigma_{NRI>0}^2$ .

$$\hat{\sigma}_{NRI>0}^2 = \frac{\hat{p}_{ne}^{up}(1 - \hat{p}_{ne}^{up})}{n_0} + \frac{\hat{p}_e^{up}(1 - \hat{p}_e^{up})}{n_1},$$

$$\sigma_{NRI>0}^2 = \frac{1}{n_0} \delta_{1,0}^2 + \frac{1}{n_1} \delta_{0,1}^2$$

From (a13),

$$\delta_{1,0}^2 = \text{Var} \left( h^{(1, 0)}(\mathbf{x}_i) \right) = \text{Var} \left( I[(\mathbf{a}^T - \mathbf{a}^* \mathbf{T}) \mathbf{x}_i > 0] \right) \approx \hat{p}_{ne}^{up}(1 - \hat{p}_{ne}^{up})$$

$$\delta_{0,1}^2 = \text{Var} \left( h^{(0, 1)}(\mathbf{y}_j) \right) = \text{Var} \left( I[(\mathbf{a}^T - \mathbf{a}^* \mathbf{T}) \mathbf{y}_j > 0] \right) \approx \hat{p}_e^{up}(1 - \hat{p}_e^{up}),$$

where  $\hat{p}_{ne}^{up} = \frac{\#non - events \text{ moving up}}{n_0}$  and  $\hat{p}_e^{up} = \frac{\#events \text{ moving up}}{n_1}$

$$\hat{\sigma}_{NRI>0}^2 = \frac{\hat{p}_{ne}^{up}(1 - \hat{p}_{ne}^{up})}{n_0} + \frac{\hat{p}_e^{up}(1 - \hat{p}_e^{up})}{n_1},$$

B.2.2.3. Derivation of  $\sigma_{3cNRI}^2$ .

$$\hat{\sigma}_{3cNRI}^2 = \frac{4(\hat{p}_{ne}^{2up} + \hat{p}_{ne}^{2down}) + \hat{p}_{ne}^{1up} + \hat{p}_{ne}^{1down} - (2(\hat{p}_{ne}^{2up} - \hat{p}_{ne}^{2down}) + \hat{p}_{ne}^{1up} - \hat{p}_{ne}^{1down})^2}{n_0} + \frac{4(\hat{p}_{ev}^{2up} + \hat{p}_{ev}^{2down}) + \hat{p}_{ev}^{1up} + \hat{p}_{ev}^{1down} - (2(\hat{p}_{ev}^{2up} - \hat{p}_{ev}^{2down}) + \hat{p}_{ev}^{1up} - \hat{p}_{ev}^{1down})^2}{n_1},$$

where  $\hat{p}_{ev}^{k\ up}$  is the fraction of events that moved  $k$  categories up in the full model.

$$\sigma_{3cNRI}^2 = \frac{1}{n_0} \delta_{1,0}^2 + \frac{1}{n_1} \delta_{0,1}^2$$

From (a15),

$$\begin{aligned} h^{(1, \ 0)}(x_i) &= -I_{js}[x_i] + const \\ h^{(0, \ 1)}(y_j) &= I_{js}[y_j] + const \end{aligned} \quad (a14)$$

$\delta_{1, \ 0}^2 = Var(h^{(1, \ 0)}(x_i)) = Var(I_{js}[x_i])$  is the variance of a multinomial random variable that takes on values in (2, 1, 0, -1, -2).

$$\begin{aligned} \hat{\delta}_{1,0}^2 &\approx 4(\hat{p}_{ne}^{2up} + \hat{p}_{ne}^{2down}) + \hat{p}_{ne}^{1up} + \hat{p}_{ne}^{1down} - (2(\hat{p}_{ne}^{2up} - \hat{p}_{ne}^{2down}) + \hat{p}_{ne}^{1up} - \hat{p}_{ne}^{1down})^2 \\ \hat{\delta}_{0,1}^2 &\approx 4(\hat{p}_{ev}^{2up} + \hat{p}_{ev}^{2down}) + \hat{p}_{ev}^{1up} + \hat{p}_{ev}^{1down} - (2(\hat{p}_{ev}^{2up} - \hat{p}_{ev}^{2down}) + \hat{p}_{ev}^{1up} - \hat{p}_{ev}^{1down})^2 \end{aligned}$$

Then

$$\sigma_{3cNRI}^2 = \frac{\hat{\delta}_{1,0}^2}{n_0} + \frac{\hat{\delta}_{0,1}^2}{n_1}, \text{ where } \hat{\delta}_{1,0}^2 \text{ and } \hat{\delta}_{0,1}^2 \text{ are defined earlier.}$$

$$\begin{aligned} \hat{\sigma}_{3cNRI}^2 &= \frac{4(\hat{p}_{ne}^{2up} + \hat{p}_{ne}^{2down}) + \hat{p}_{ne}^{1up} + \hat{p}_{ne}^{1down} - (2(\hat{p}_{ne}^{2up} - \hat{p}_{ne}^{2down}) + \hat{p}_{ne}^{1up} - \hat{p}_{ne}^{1down})^2}{n_0} \\ &+ \frac{4(\hat{p}_{ev}^{2up} + \hat{p}_{ev}^{2down}) + \hat{p}_{ev}^{1up} + \hat{p}_{ev}^{1down} - (2(\hat{p}_{ev}^{2up} - \hat{p}_{ev}^{2down}) + \hat{p}_{ev}^{1up} - \hat{p}_{ev}^{1down})^2}{n_1} \end{aligned}$$

#### B.2.2.4. Derivation of $\sigma_{IDI}^2$ .

$$\begin{aligned} \hat{\sigma}_{IDI}^2 &= \frac{Var(\Delta predp(x_i))}{n_0} + \frac{Var(\Delta predp(y_j))}{n_1}, \\ \sigma_{IDI}^2 &= \frac{1}{n_0} \delta_{1,0}^2 + \frac{1}{n_1} \delta_{0,1}^2 \end{aligned}$$

From (a14),

$$\begin{aligned} \delta_{1,0}^2 &= Var(h^{(1, \ 0)}(x_i)) = Var(f(a' x_i) - f(a^* x_i)) \approx Var(\Delta predp(x_i)) \\ \delta_{0,1}^2 &= Var(h^{(0, \ 1)}(y_j)) = Var(f(a' y_j) - f(a^* y_j)) \approx Var(\Delta predp(y_j)) \end{aligned}$$

where

$\Delta predp(x_i)$  is the difference in predicted probabilities from two models for non-event  $x_i$  and  $\Delta predp(y_j)$  is the difference in predicted probabilities from two models for event  $y_j$

$$\hat{\sigma}_{IDI}^2 = \frac{Var(\Delta predp(x_i))}{n_0} + \frac{Var(\Delta predp(y_j))}{n_1},$$

Variances derived here are summarized in Table III of the paper and are reproduced here (Table A1). q.e.d.

**Table A1.** Variance formulas in non-degenerate case, unadjusted for estimated parameters.

	$\hat{\sigma}^2$ , ignoring the adjustment for estimated parameters	Requires adjustment?
$\hat{\sigma}_{\Delta AUC}^2$ no tied ranks	$\frac{\text{Var}(\text{rank}_e^*(a^T \mathbf{x}_i) - \text{rank}_e(a^T \mathbf{x}_i))}{n_0} + \frac{\text{Var}(\text{rank}_{ne}^*(a^T \mathbf{y}_j) - \text{rank}_{ne}(a^T \mathbf{y}_j))}{n_1}$	No
$\hat{\sigma}_{\Delta AUC}^2$ tied ranks	Use DeLong formula [8]	No
$\hat{\sigma}_{NRI>0}^2$	$\frac{\hat{p}_{ne}^{up}(1 - \hat{p}_{ne}^{up})}{n_0} + \frac{\hat{p}_e^{up}(1 - \hat{p}_e^{up})}{n_1}$	No
$\hat{\sigma}_{NRI(r)}^2$	$\frac{n_0}{\hat{p}_{ne}^{up} + \hat{p}_{ne}^{down} - (\hat{p}_{ne}^{up} - \hat{p}_{ne}^{down})^2} + \frac{n_1}{\hat{p}_{ev}^{up} + \hat{p}_{ev}^{down} - (\hat{p}_{ev}^{up} - \hat{p}_{ev}^{down})^2}$	No
$\hat{\sigma}_{3cNRI}^2$	$\frac{4(\hat{p}_{ne}^{2up} + \hat{p}_{ne}^{2down}) + \hat{p}_{ne}^{1up} + \hat{p}_{ne}^{1down} - (2(\hat{p}_{ne}^{2up} - \hat{p}_{ne}^{2down}) + \hat{p}_{ne}^{1up} - \hat{p}_{ne}^{1down})^2}{n_0} + \frac{4(\hat{p}_{ev}^{2up} + \hat{p}_{ev}^{2down}) + \hat{p}_{ev}^{1up} + \hat{p}_{ev}^{1down} - (2(\hat{p}_{ev}^{2up} - \hat{p}_{ev}^{2down}) + \hat{p}_{ev}^{1up} - \hat{p}_{ev}^{1down})^2}{n_1}$	Yes
$\hat{\sigma}_{IDI}^2$	$\frac{\text{Var}(\Delta \text{pred}p(\mathbf{x}_i))}{n_0} + \frac{\text{Var}(\Delta \text{pred}p(\mathbf{y}_j))}{n_1}$	Yes

## Appendix C

### C.1. Closed-form expression for 3cNRI under the assumption of normality of predictor variables

Closed-form expression of 3cNRI under the assumption of normality of predictor variables was derived in [9].

$$\begin{aligned}
 3cNRI = & \Phi\left(\frac{\frac{M_{new}^2}{2} - \ln\left(\frac{t_H(1-r)}{(1-t_H)r}\right)}{\sqrt{M_{new}^2}}\right) - \Phi\left(\frac{\frac{M_{old}^2}{2} - \ln\left(\frac{t_H(1-r)}{(1-t_H)r}\right)}{\sqrt{M_{old}^2}}\right) + \Phi\left(\frac{\frac{M_{new}^2}{2} - \ln\left(\frac{t_L(1-r)}{(1-t_L)r}\right)}{\sqrt{M_{new}^2}}\right) \\
 & - \Phi\left(\frac{\frac{M_{old}^2}{2} - \ln\left(\frac{t_L(1-r)}{(1-t_L)r}\right)}{\sqrt{M_{old}^2}}\right) + \Phi\left(\frac{\frac{M_{new}^2}{2} + \ln\left(\frac{t_H(1-r)}{(1-t_H)r}\right)}{\sqrt{M_{new}^2}}\right) - \Phi\left(\frac{\frac{M_{old}^2}{2} + \ln\left(\frac{t_H(1-r)}{(1-t_H)r}\right)}{\sqrt{M_{old}^2}}\right) \\
 & + \Phi\left(\frac{\frac{M_{new}^2}{2} + \ln\left(\frac{t_L(1-r)}{(1-t_L)r}\right)}{\sqrt{M_{new}^2}}\right) - \Phi\left(\frac{\frac{M_{old}^2}{2} + \ln\left(\frac{t_L(1-r)}{(1-t_L)r}\right)}{\sqrt{M_{old}^2}}\right)
 \end{aligned}
 \tag{A3.1}$$

where  $M_{new}^2$ ,  $M_{old}^2$  are squared Mahalanobis distance for new and old models,  $t_H$ ,  $t_L$ ,  $r$  are high and low thresholds and an event rate.

### C.2. Closed-form expression for IDI under the assumption of normality of predictor variables

Closed-form expression of IDI under the assumption of normality of predictor variables was derived in [39].

$$\begin{aligned}
 IDI = & \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi M_{new}^2}} \exp\left(-\frac{(x - 0.5M_{new}^2)^2}{2M_{new}^2}\right) \left(\frac{1}{1 + \frac{r}{1-r} \exp(-x)} - \frac{1}{1 + \frac{r}{1-r} \exp(x)}\right) dx \\
 & - \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi M_{old}^2}} \exp\left(-\frac{(x - 0.5M_{old}^2)^2}{2M_{old}^2}\right) \left(\frac{1}{1 + \frac{r}{1-r} \exp(-x)} - \frac{1}{1 + \frac{r}{1-r} \exp(x)}\right) dx,
 \end{aligned}
 \tag{A3.2}$$

where  $r$  is the event rate.

## Appendix D

Derivative of closed-form formula of 3cNRI [formula (A3.1)] with respect to event rate  $r$  is as follows:

$$(3cNRI)'_r = \frac{1}{\sqrt{2\pi}(r-1)r} \times \left( \frac{1}{M_{new}} \left( \exp \left( -\frac{\left( M_{new}^2 - 2 \ln \left[ \frac{(1-r)t_H}{r(1-t_H)} \right] \right)^2}{8M_{new}^2} \right) - \exp \left( -\frac{\left( M_{new}^2 + 2 \ln \left[ \frac{(1-r)t_H}{r(1-t_H)} \right] \right)^2}{8M_{new}^2} \right) + \right. \right. \\ \left. \exp \left( -\frac{\left( M_{new}^2 - 2 \ln \left[ \frac{(1-r)t_L}{r(1-t_L)} \right] \right)^2}{8M_{new}^2} \right) - \exp \left( -\frac{\left( M_{new}^2 + 2 \ln \left[ \frac{(1-r)t_L}{r(1-t_L)} \right] \right)^2}{8M_{new}^2} \right) \right) \\ \left. + \frac{1}{M_{old}} \left( -\exp \left( -\frac{\left( M_{old}^2 - 2 \ln \left[ \frac{(1-r)t_H}{r(1-t_H)} \right] \right)^2}{8M_{old}^2} \right) + \exp \left( -\frac{\left( M_{old}^2 - 2 \ln \left[ \frac{(1-r)t_H}{r(1-t_H)} \right] \right)^2}{8M_{old}^2} \right) - \right. \right. \\ \left. \exp \left( -\frac{\left( M_{old}^2 - 2 \ln \left[ \frac{(1-r)t_L}{r(1-t_L)} \right] \right)^2}{8M_{old}^2} \right) + \exp \left( -\frac{\left( M_{old}^2 + 2 \ln \left[ \frac{(1-r)t_L}{r(1-t_L)} \right] \right)^2}{8M_{old}^2} \right) \right) \right)$$

Derivative of closed-form formula of IDI [formula (A3.2)] with respect to event rate  $r$  is as follows:

$$(IDI)'_r = \frac{1}{\sqrt{2\pi}M_{new}M_{old}} \times \int_{-\infty}^{+\infty} \frac{(-1 + e^{2x})(2r-1)}{(r + e^x(1-r))^2(1 + (e^x - 1)r)^2} \exp \left( x - \frac{(x - 0.5M_{old}^2)^2}{2M_{old}^2} - \frac{(x - 0.5M_{new}^2)^2}{2M_{new}^2} \right) \\ \left( M_{new} \exp \left( \frac{(x - 0.5M_{new}^2)^2}{2M_{new}^2} - \frac{(x - 0.5M_{new}^2)^2}{2M_{new}^2} \right) - M_{old} \exp \left( \frac{(x - 0.5M_{old}^2)^2}{2M_{old}^2} - \frac{(x - 0.5M_{new}^2)^2}{2M_{new}^2} \right) \right) dx$$

1. Lee J. U-Statistics: Theory and Practice. 1990.
2. Hoeffding W. A class of statistics with asymptotically normal distribution. *The annals of mathematical statistics* 1948;293–325.
3. Pencina MJ, D'Agostino RB, and Vasan RS. Evaluating the added predictive ability of a new marker: from area under the roc curve to reclassification and beyond. *Stat Med* 2008; **2**:157–172.
4. Pencina MJ, D'Agostino RB, and Steyerberg EW. Extensions of net reclassification improvement calculations to measure usefulness of new biomarkers. *Stat Med* 2011; **1**:11–21.
5. Pepe MS. *The Statistical Evaluation of Medical Tests for Classification and Prediction*. Oxford University Press, USA, 2003.
6. Bamber D. The area above the ordinal dominance graph and the area below the receiver operating characteristic graph. *Journal of mathematical psychology* 1975; **4**:387–415.
7. Hanley JA, and McNeil BJ. The meaning and use of the area under a receiver operating characteristic (ROC) curve. *Radiology* 1982; **1**:29–36.
8. DeLong ER, DeLong DM, and Clarke-Pearson DL. Comparing the areas under two or more correlated receiver operating characteristic curves: a nonparametric approach. *Biometrics* 1988;837–845.
9. Pencina KM, Pencina MJ, and D'Agostino RB. What to expect from net reclassification improvement with three categories. *Stat Med* 2014; **28**:4975–4987.

## References

1. Wilson PW, D'Agostino RB, Levy D, Belanger AM, Silbershatz H, Kannel WB. Prediction of coronary heart disease using risk factor categories. *Circulation* 1998; **18**:1837–1847.
2. Expert Panel on Detection E. Executive summary of the third report of the National Cholesterol Education Program (NCEP) expert panel on detection, evaluation, and treatment of high blood cholesterol in adults (adult treatment panel III). *JAMA* 2001; **19**:2486.
3. Gail MH, Brinton LA, Byar DP, Corle DK, Green SB, Schairer C, Mulvihill JJ. Projecting individualized probabilities of developing breast cancer for white females who are being examined annually. *Journal of the National Cancer Institute* 1989; **24**:1879–1886.
4. Dehmer SP, Maciosek MV, Flottemesch TJ. Aspirin use to prevent cardiovascular disease and colorectal cancer. *Evidence Syntheses* 2015; **2015**:131s.
5. Stone NJ, Robinson J, Lichtenstein AH, Bairey Merz CN, Blum CB, Eckel RH, Goldberg AC, Gordon D, Levy D, Lloyd-Jones DM, McBride P, Schwartz JS, Shero ST, Smith SC Jr, Watson K, Wilson PWF. 2013 ACC/AHA guideline on the treatment of blood cholesterol to reduce atherosclerotic cardiovascular risk in adults: a report of the american college of cardiology/american heart association task force on practice guidelines. *Journal of the American College of Cardiology* 2014; **25**\_PA:2889–2934.
6. Bamber D. The area above the ordinal dominance graph and the area below the receiver operating characteristic graph. *Journal of Mathematical Psychology* 1975; **4**:387–415.
7. Hanley JA, McNeil BJ. The meaning and use of the area under a receiver operating characteristic (ROC) curve. *Radiology* 1982; **1**:29–36.
8. Pencina MJ, D'Agostino RB, Vasan RS. Evaluating the added predictive ability of a new marker: from area under the ROC curve to reclassification and beyond. *Statistics in Medicine* 2008; **2**:157–172.
9. Pencina MJ, D'Agostino RB, Steyerberg EW. Extensions of net reclassification improvement calculations to measure usefulness of new biomarkers. *Statistics in Medicine* 2011; **1**:11–21.
10. Pencina MJ, Neely B, Steyerberg EW. Re: net risk reclassification P values: valid or misleading? *Journal of the National Cancer Institute* 2015; **1** dju355.
11. DeLong ER, DeLong DM, Clarke-Pearson DL. Comparing the areas under two or more correlated receiver operating characteristic curves: a nonparametric approach. *Biometrics* 1988; **44**:837–845.
12. Demler OV, Pencina MJ, D'Agostino RB Sr. Misuse of DeLong test to compare AUCs for nested models. *Statistics in Medicine* 2012; **23**:2577–2587.
13. Kerr KF, Wang Z, Janes H, McClelland RL, Psaty BM, Pepe MS. Net reclassification indices for evaluating risk-prediction instruments: a critical review. *Epidemiology (Cambridge, Mass.)* 2014; **1**:114.
14. Seshan VE, Gönen M, Begg CB. Comparing ROC curves derived from regression models. *Statistics in Medicine* 2013; **9**:1483–1493.
15. Pencina MJ, Steyerberg EW, D'Agostino RB. Net reclassification index at event rate: properties and relationships. *Statistics in Medicine* 2017. DOI: 10.1002/sim.7041
16. Pencina KM, Pencina MJ, D'Agostino RB. What to expect from net reclassification improvement with three categories. *Statistics in Medicine* 2014; **28**:4975–4987.
17. Klein JP. Small sample moments of some estimators of the variance of the Kaplan–Meier and Nelson–Aalen estimators. *Scandinavian Journal of Statistics* 1991; **18**:333–340.
18. Antolini L, Nam B-H, D'Agostino RB. Inference on correlated discrimination measures in survival analysis: a nonparametric approach. *Communications in statistics-Theory and Methods* 2004; **9**:2117–2135.
19. Lee MLT, Rosner BA. The average area under correlated receiver operating characteristic curves: a nonparametric approach based on generalized two-sample Wilcoxon statistics. *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 2001; **3**:337–344.
20. Kerr KF, McClelland RL, Brown ER, Lumley T. Evaluating the incremental value of new biomarkers with integrated discrimination improvement. *American Journal of Epidemiology* 2011; **3**:364–374.
21. Harrell FE, Lee KL, Califf RM, Pryor DB, Rosati RA. Regression modelling strategies for improved prognostic prediction. *Statistics in Medicine* 1984; **2**:143–152.
22. Pepe MS. The Statistical Evaluation of Medical Tests for Classification and Prediction. Oxford University Press: USA, 2003.
23. Pepe M, Janes H. Methods for evaluating prediction performance of biomarkers and tests. In Risk Assessment and Evaluation of Predictions. Springer Science+Business Media: New York, 2013; 107–142.
24. So H-C, Sham PC. A unifying framework for evaluating the predictive power of genetic variants based on the level of heritability explained. *PLoS Genetics* 2010; **12** e1001230.
25. Pencina MJ, Fine JP, D'Agostino RB. Discrimination slope and integrated discrimination improvement–properties, relationships and impact of calibration. *Statistics in Medicine* 2016 in press.
26. Pepe MS, Kerr KF, Longton G, Wang Z. Testing for improvement in prediction model performance. *Statistics in Medicine* 2013; **9**:1467–1482.
27. Demler OV, Pencina MJ, D'Agostino RB Sr. Equivalence of improvement in area under ROC curve and linear discriminant analysis coefficient under assumption of normality. *Statistics in Medicine* 2011; **12**:1410–1418.
28. Lee J. U-statistics: theory and practice. Marcel Dekker: New York, 1990.
29. Hoeffding W. A class of statistics with asymptotically normal distribution. *The Annals of Mathematical Statistics* 1948; **19**:293–325.
30. Lehmann EL. Consistency and unbiasedness of certain nonparametric tests. *The Annals of Mathematical Statistics* 1951; **22**:165–179.
31. Serfling RJ. Approximation Theorems of Mathematical Statistics, Vol. **162**. John Wiley & Sons: New York, 2009.



32. Sukhatme BV. Testing the hypothesis that two populations differ only in location. *The Annals of Mathematical Statistics* 1958; **29**:60–78.
33. Randles RH. On the asymptotic normality of statistics with estimated parameters. *The Annals of Statistics* 1982; **10**:462–474.
34. de Wet T, Randles RH. On the effect of substituting parameter estimators in limiting  $X^2$  U and V statistics. *The Annals of Statistics* 1987; **15**:398–412.
35. Efron B. The efficiency of logistic regression compared to normal discriminant analysis. *Journal of the American Statistical Association* 1975; **352**:892–898.
36. Su JQ, Liu JS. Linear combinations of multiple diagnostic markers. *Journal of the American Statistical Association* 1993; **424**:1350–1355.
37. Mardia KV, Kent JT, Bibby JM. *Multivariate Analysis (Probability and Mathematical Statistics)*. Academic Press: London, 1980.
38. Pencina MJ, D'Agostino RB Sr, Demler OV. Novel metrics for evaluating improvement in discrimination: net reclassification and integrated discrimination improvement for normal variables and nested models. *Statistics in Medicine* 2012; **2**:101–113.
39. Paynter NP, Cook NR. A bias-corrected net reclassification improvement for clinical subgroups. *Medical Decision Making* 2013; **2**:154–162.
40. D'Agostino RB, Vasan RS, Pencina MJ, Wolf PA, Cobain M, Massaro JM, Kannel WB. General cardiovascular risk profile for use in primary care the framingham heart study. *Circulation* 2008; **6**:743–753.
41. Van Der Vaart AW, Wellner JA. Weak Convergence. In *Weak Convergence and Empirical Processes*. Springer: New York, 1996; 16–28.
42. Babu GJ, Singh K. On one term Edgeworth correction by Efron's bootstrap. *Sankhyā: The Indian Journal of Statistics, Series A* 1984; **46**:219–232.
43. Singh K. On the asymptotic accuracy of Efron's bootstrap. *The Annals of Statistics* 1981; **9**:1187–1195.