
COMPARISON OF THE PREDICTIVE VALUES OF MULTIPLE BINARY DIAGNOSTIC TESTS IN THE PRESENCE OF IGNORABLE MISSING DATA

Authors: ANA EUGENIA MARÍN-JIMÉNEZ
– Statistics (Biostatistics), Faculty of Education, Economy and Technology,
University of Granada, Ceuta, Spain
anamarin@ugr.es

JOSÉ ANTONIO ROLDÁN-NOFUENTES
– Statistics (Biostatistics), School of Medicine, University of Granada,
Granada, Spain
jaroldan@ugr.es

Received: August 2013

Revised: April 2015

Accepted: May 2015

Abstract:

- The comparison of the predictive values of binary diagnostic tests is an important topic in the study of statistical methods applied to medical diagnosis. In this article, we study a global hypothesis test to simultaneously compare the predictive values of multiple binary diagnostic tests in the presence of ignorable missing data. The global hypothesis test deduced is based on the chi-squared distribution. Simulation experiments were carried out to study the type I error probability and the power of global hypothesis test and of other alternative methods when comparing the predictive values of two and three binary diagnostic tests respectively.

Key-Words:

- *global hypothesis test; predictive values; multiple comparisons; chi-squared distribution; ignorable missing data.*

AMS Subject Classification:

- 62P10, 62H15.

1. INTRODUCTION

A diagnostic test is a medical test that is applied to an individual in order to determine the presence or absence of a disease. When the diagnostic test can only give two possible results (positive, indicating the provisional presence of the disease, or negative, indicating the provisional absence of the disease) the diagnostic test is called a binary diagnostic test (*BDT*) and it is used very frequently in clinical practice. A stress test for the diagnosis of coronary disease or a mammogram to diagnose breast cancer are two examples of *BDTs*. The most common parameters to assess the accuracy or performance of a *BDT* are sensitivity (*Se*) and specificity (*Sp*). Other commonly used parameters to assess the performance of a *BDT* are predictive values (*PVs*). The positive predictive value (*PPV*) of a *BDT* is the probability of an individual having the disease given that the result of the *BDT* is positive and the negative predictive value (*NPV*) is the probability of an individual not having the disease given that the result of the *BDT* is negative. The predictive values (*PVs*) are a measure of clinical accuracy of the *BDT*, and they depend on the sensitivity and the specificity of the *BDT* and on the prevalence of the disease (*p*). Applying Bayes Theorem, the *PVs* are calculated as

$$(1.1) \quad \begin{aligned} &PPV = \frac{p \times Se}{p \times Se + (1-p) \times (1-Sp)} \\ &\text{and} \\ &NPV = \frac{(1-p) \times Sp}{p \times (1-Se) + (1-p) \times Sp} . \end{aligned}$$

In the study of statistical methods for the diagnosis of diseases, comparison of the accuracy or the performance of two diagnostic tests is a topic of particular importance. In paired designs (*i.e.* when the two *BDTs* and the gold standard are applied to all of the individuals in a random sample), comparison of the *PVs* of two *BDTs* in relation to the same gold standard has been the subject of several studies in the statistical literature [1, 2, 3, 4]. In all of them, the comparison of the two *PPVs* and the comparison of the two *NPVs* is carried out independently. Roldán-Nofuentes *et al.* [5] showed that the *PVs* of two (or more) *BDTs* are correlated and they studied a global hypothesis test based on the chi-squared distribution to simultaneously compare the *PVs* of two or more *BDTs* in relation to the same gold standard. In all of these studies, the disease status of all of the patients is known, as well as the results of the two diagnostic tests. This situation is also known as ‘complete verification’ (because the gold standard is applied to all of the individuals in the sample). Poleto *et al.* [6] studied the comparison of the predictive values of two *BDTs* when for some individuals we do not know the results of one of the two *BDTs*. Furthermore, in clinical practice it is common for the gold standard not to be applied to all of the individuals in the sample, thus leading to the problem known as partial disease verification [7, 8, 9].

Therefore, the disease status (if the disease is present or absent) is unknown for a subset of individuals in the sample. In this situation, Roldán-Nofuentes *et al.* [10, 11] studied the comparison of the *PPVs* and of the *NPVs* of two *BDTs*. Nevertheless, in these two studies they did not consider the dependence that exists between the *PVs* of the diagnostic tests. This is the essence of our article, to study a global hypothesis test that allows us to jointly compare the *PVs* of two (or more) *BDTs* in the presence of ignorable missing data. In this article, we study a global hypothesis test to simultaneously compare the predictive values of two or more *BDTs* when, in the presence of partial disease verification, the missing data mechanism is ignorable. In Section 2, we propose a global hypothesis test, and other alternative methods, to simultaneously compare the *PVs* of multiple *BDTs*. In Section 3, Monte Carlo simulation experiments are carried out in order to study the type I error probability and the power of the global hypothesis test (and of the alternative methods) when comparing the *PVs* of two and of three *BDTs* respectively. In Section 4, the method proposed is applied to two examples, and in Section 5 the results obtained are discussed.

2. THE MODEL

Let us consider J *BDTs* ($J \geq 2$) that are applied independently to the same random sample of size n extracted from a population that has a determined prevalence of the disease (p). Moreover, let us consider that the gold standard has not been applied to all of the individuals in the random sample. In this situation, the J diagnostic tests are applied to all of the individuals in the sample whilst the gold standard is only applied to a subset of them. Therefore, the results of the J diagnostic tests are known by all of the individuals in the sample, whereas the result of the gold standard (*i.e.* the disease status) is only unknown to a subset of them. Let T_j , V and D be the random binary variables defined as: T_j which models the result of the j -th *BDT* ($j = 1, \dots, J$), so that $T_j = 1$ when the test result is positive and $T_j = 0$ when the result is negative; V models the verification process, $V = 1$ when the individual is verified with the gold standard and $V = 0$ when the individual is not verified; and D models the result of the gold standard, $D = 1$ when the individual has the disease and $D = 0$ when the individual does not. Let $Se_j = P(T_j = 1 | D = 1)$, $Sp_j = P(T_j = 0 | D = 0)$, $PPV_j = P(D = 1 | T_j = 1)$ and $NPV_j = P(D = 0 | T_j = 0)$ be the sensitivity, the specificity, the positive predictive value and the negative predictive value of the j -th *BDT* respectively. Let the observed frequencies be: s_{i_1, \dots, i_J} is the number of patients verified in which $T_1 = i_1, T_2 = i_2, \dots, T_J = i_J$ and $D = 1$; r_{i_1, \dots, i_J} is the number of patients verified in which $T_1 = i_1, T_2 = i_2, \dots, T_J = i_J$ and $D = 0$; and u_{i_1, \dots, i_J} is the number of patients not verified in which $T_1 = i_1, T_2 = i_2, \dots, T_J = i_J$ with $i_j = 0, 1$ and $j = 1, \dots, J$. Let $n_{i_1, \dots, i_J} = s_{i_1, \dots, i_J} + r_{i_1, \dots, i_J} + u_{i_1, \dots, i_J}$ and

$n = \sum_{i_1, \dots, i_J=0}^1 n_{i_1, \dots, i_J}$. As only a subset of individuals in the sample have their disease status verified with the gold standard, the verification probabilities ($\lambda_{k, i_1, \dots, i_J}$) are defined as the probability of selecting an individual for whom $D = k$, $T_1 = i_1$, $T_2 = i_2$, ..., $T_J = i_J$ with $k, i_j = 0, 1$, $j = 1, \dots, J$, to verify his or her disease status i.e.

$$\lambda_{k, i_1, \dots, i_J} = P\left(V = 1 \mid D = k, T_1 = i_1, T_2 = i_2, \dots, T_J = i_J\right).$$

Assuming that the verification process with the gold standard only depends on the results of the J BDTs and does not depend on the disease status, then the missing data mechanism is missing at random (MAR) [12]. Assuming also that the parameters of the data model and the parameters of the missingness mechanism are distinct, then the missing data mechanism is ignorable [13]. Under this model, the verification probabilities are

$$\lambda_{k, i_1, \dots, i_J} = \lambda_{i_1, \dots, i_J} = P\left(V = 1 \mid T_1 = i_1, T_2 = i_2, \dots, T_J = i_J\right),$$

and all of the parameters can be estimated applying the maximum likelihood method.

2.1. Maximum likelihood estimators of the PVs

As the J BDTs are applied to all of the n individuals in the random sample and the gold standard is only applied to a subset of them, the frequencies observed r_{i_1, \dots, i_J} , s_{i_1, \dots, i_J} and u_{i_1, \dots, i_J} with $i_j = 0, 1$ and $j = 1, \dots, J$, which can be written in the form of a 3×2^J table in which the sample of size n has been set, are the realization of a multinomial distribution whose probabilities are

$$\phi_{i_1, \dots, i_J} = P\left(V = 1, D = 1, T_1 = i_1, T_2 = i_2, \dots, T_J = i_J\right),$$

$$\varphi_{i_1, \dots, i_J} = P\left(V = 1, D = 0, T_1 = i_1, T_2 = i_2, \dots, T_J = i_J\right)$$

and

$$\gamma_{i_1, \dots, i_J} = P\left(V = 0, T_1 = i_1, T_2 = i_2, \dots, T_J = i_J\right).$$

Let $\boldsymbol{\omega} = (\phi_{1, \dots, 1}, \dots, \phi_{0, \dots, 0}, \varphi_{1, \dots, 1}, \dots, \varphi_{0, \dots, 0}, \gamma_{1, \dots, 1}, \dots, \gamma_{0, \dots, 0})^T$ be a vector sized $(3 \cdot 2^J)$ whose components are the probabilities of multinomial distribution and $\eta_{i_1, \dots, i_J} = \phi_{i_1, \dots, i_J} + \varphi_{i_1, \dots, i_J} + \gamma_{i_1, \dots, i_J}$. Assuming that the missing data mechanism is ignorable, the PVs of the j -th BDT are written in terms of the parameters of

the vector ω and of the verification probabilities as

$$PPV_j = \frac{\sum_{i_1, \dots, i_J=0; i_j=1}^1 \phi_{i_1, \dots, i_J} \lambda_{i_1, \dots, i_J}^{-1}}{\sum_{i_1, \dots, i_J=0; i_j=1}^1 \eta_{i_1, \dots, i_J}} \quad (2.1)$$

and

$$NPV_j = \frac{\sum_{i_1, \dots, i_J=0; i_j=0}^1 \varphi_{i_1, \dots, i_J} \lambda_{i_1, \dots, i_J}^{-1}}{\sum_{i_1, \dots, i_J=0; i_j=0}^1 \eta_{i_1, \dots, i_J}},$$

where $\lambda_{i_1, \dots, i_J} = (\phi_{i_1, \dots, i_J} + \varphi_{i_1, \dots, i_J}) / \eta_{i_1, \dots, i_J}$ are the verification probabilities. Therefore, in equations (2.1) we can observe the dependence of the *PVs* of the verification process subject to the MAR assumption. In this situation the logarithm of the likelihood function is

$$l = \sum_{i_1, \dots, i_J=0}^1 s_{i_1, \dots, i_J} \log(\phi_{i_1, \dots, i_J}) + \sum_{i_1, \dots, i_J=0}^1 r_{i_1, \dots, i_J} \log(\varphi_{i_1, \dots, i_J}) \\ + \sum_{i_1, \dots, i_J=0}^1 u_{i_1, \dots, i_J} \log(\gamma_{i_1, \dots, i_J}),$$

so that maximizing this function, the maximum likelihood estimators (MLEs) of ϕ_{i_1, \dots, i_J} , $\varphi_{i_1, \dots, i_J}$ and γ_{i_1, \dots, i_J} are the estimators of multinomial proportions [14], i.e.

$$(2.2) \quad \hat{\phi}_{i_1, \dots, i_J} = \frac{s_{i_1, \dots, i_J}}{n}, \quad \hat{\varphi}_{i_1, \dots, i_J} = \frac{r_{i_1, \dots, i_J}}{n} \quad \text{and} \quad \hat{\gamma}_{i_1, \dots, i_J} = \frac{u_{i_1, \dots, i_J}}{n},$$

and the *MLE* of η_{i_1, \dots, i_J} is $\hat{\eta}_{i_1, \dots, i_J} = n_{i_1, \dots, i_J} / n$. Substituting in equations (2.1) the parameters with their respective *MLEs* given in equations (2.2), the *MLEs* of the *PVs* of the *j*-th *BDT* are

$$\widehat{PPV}_j = \frac{\sum_{i_1, \dots, i_J=0; i_j=1}^1 \frac{s_{i_1, \dots, i_J} n_{i_1, \dots, i_J}}{s_{i_1, \dots, i_J} + r_{i_1, \dots, i_J}}}{\sum_{i_1, \dots, i_J=0; i_j=1}^1 n_{i_1, \dots, i_J}}$$

and

$$\widehat{NPV}_j = \frac{\sum_{i_1, \dots, i_J=0; i_j=0}^1 \frac{r_{i_1, \dots, i_J} n_{i_1, \dots, i_J}}{s_{i_1, \dots, i_J} + r_{i_1, \dots, i_J}}}{\sum_{i_1, \dots, i_J=0; i_j=0}^1 n_{i_1, \dots, i_J}}.$$

Once we have obtained the *MLEs* of the *PVs* of the *J BDTs*, we then estimate their variances-covariances.

2.2. Estimation of the variances-covariances of the PVs

As the vector $\boldsymbol{\omega}$ is the vector of probabilities of a multinomial distribution, the variance-covariance matrix of $\hat{\boldsymbol{\omega}}$ is $\sum_{\hat{\boldsymbol{\omega}}} = \{\text{diag}(\boldsymbol{\omega}) - \boldsymbol{\omega}^T \boldsymbol{\omega}\} / n$. Let $\boldsymbol{\tau} = (PPV_1, \dots, PPV_J, NPV_1, \dots, NPV_J)^T$ be a vector sized $2J$ whose components are the PVs of the J BDTs, and let $\hat{\boldsymbol{\tau}}$ be the MLE of $\boldsymbol{\tau}$. As $\boldsymbol{\tau}$ is a function of the components of the vector $\boldsymbol{\omega}$, applying the delta method [15] the asymptotic variance-covariance matrix of $\hat{\boldsymbol{\tau}}$ is

$$\sum_{\hat{\boldsymbol{\tau}}} = \left(\frac{\partial \boldsymbol{\tau}}{\partial \boldsymbol{\omega}} \right) \sum_{\hat{\boldsymbol{\omega}}} \left(\frac{\partial \boldsymbol{\tau}}{\partial \boldsymbol{\omega}} \right)^T.$$

Substituting in the previous expression each parameter with its corresponding MLE, we obtain the estimated asymptotic variances-covariances of the estimators of the PVs of the J BDTs.

Moreover, the asymptotic variances-covariances of $\hat{\boldsymbol{\tau}}$ can also be estimated through bootstrap [16], generating, from the random sample of size n , B samples with replacement and from these B samples asymptotic variance-covariance matrix of $\hat{\boldsymbol{\tau}}$ is estimated.

Once we have obtained the MLEs of the PVs and their estimated asymptotic variances-covariances, it is possible to solve the global hypothesis test to simultaneously compare the PVs of the J BDTs.

2.3. Global hypothesis test

The global hypothesis test to simultaneously compare the PVs of J BDTs is

$$\begin{aligned} H_0: & PPV_1 = PPV_2 = \dots = PPV_J \quad \text{and} \quad NPV_1 = NPV_2 = \dots = NPV_J, \\ H_1: & \text{at least one equality is not true,} \end{aligned}$$

which is equivalent to the hypothesis test

$$(2.3) \quad H_0: \mathbf{A}\boldsymbol{\tau} = 0 \quad \text{vs} \quad H_1: \mathbf{A}\boldsymbol{\tau} \neq 0,$$

where \mathbf{A} is a full rank matrix sized $2(J-1) \times 2J$ whose elements are known constants. For two BDTs ($J=2$) the matrix \mathbf{A} is $\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \otimes (1 \ -1)$, and for $J=3$ this matrix is $\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \otimes \begin{pmatrix} 1 & -1 & 0 \\ 0 & 1 & -1 \end{pmatrix}$, where \otimes is the Kronecker product. As the vector $\hat{\boldsymbol{\tau}}$ is asymptotically distributed according to a normal multivariate

distribution, i.e. $\hat{\boldsymbol{\tau}} \xrightarrow[n \rightarrow \infty]{} N(\boldsymbol{\tau}, \boldsymbol{\Sigma}_{\boldsymbol{\tau}})$, the Wald statistic for the global hypothesis test (2.3) is

$$(2.4) \quad Q^2 = \hat{\boldsymbol{\tau}}^T \mathbf{A}^T \left(\mathbf{A} \hat{\boldsymbol{\Sigma}}_{\hat{\boldsymbol{\tau}}} \mathbf{A}^T \right)^{-1} \mathbf{A} \hat{\boldsymbol{\tau}},$$

which is asymptotically distributed according to a chi-squared distribution with $2(J - 1)$ degrees of freedom when the null hypothesis is true.

If the global hypothesis test is solved applying bootstrap, the statistic for the global test is similar to that given in the expression (2.4), substituting $\hat{\boldsymbol{\tau}}$ with the bootstrap estimator of $\boldsymbol{\tau}$ and $\hat{\boldsymbol{\Sigma}}_{\hat{\boldsymbol{\tau}}}$ with the variance-covariance matrix estimated through bootstrap.

Other alternative methods will now be proposed to solve the global hypothesis test (2.3).

2.4. Alternative methods

The method proposed in the previous Section to solve the global hypothesis test (2.3) is based on the chi-squared distribution. The following are some alternative methods to solve this hypothesis test:

Method 1. Consists of solving the $J(J - 1)$ marginal hypothesis tests given by

$$H_0: PPV_k = PPV_l \quad \text{vs} \quad H_1: PPV_k \neq PPV_l$$

and

$$H_0: NPV_k = NPV_l \quad \text{vs} \quad H_1: NPV_k \neq NPV_l$$

with $k, l = 1, \dots, J$ and $k \neq l$, each one to an error rate of $\alpha/\{J(J - 1)\}$, i.e. applying the Bonferroni method [17], where the statistics is

$$(2.5) \quad z = \frac{\widehat{PV}_k - \widehat{PV}_l}{\sqrt{\widehat{\text{Var}}(\widehat{PV}_k) + \widehat{\text{Var}}(\widehat{PV}_l) - 2\widehat{\text{Cov}}(\widehat{PV}_k, \widehat{PV}_l)}} \longrightarrow N(0, 1),$$

and where PV is PPV or NPV respectively.

Method 2. Consists of solving the $J(J - 1)$ marginal hypothesis tests and applying the multiple comparison method of Holm [18] to a global error rate of α .

Method 3. Consists of solving the $J(J - 1)$ marginal hypothesis tests and applying the multiple comparison method of Hochberg [19] to a global error rate of α .

These methods are very easy to apply from the p-values calculated in the $J(J - 1)$ marginal hypothesis tests. The Bonferroni method is a classic method of post hoc comparison, and the Holm and Hochberg methods are less conservative post hoc methods than the Bonferroni method. Furthermore all of the aforementioned methods can be applied both if the *PVs* and their variances-covariances are estimated through the maximum likelihood method and the delta method respectively, or if they are estimated through the bootstrap method.

3. SIMULATION EXPERIMENTS

Monte Carlo simulation experiments were carried out to study the type I error probability and the power of the global hypothesis proposed in Section 2.3 and of the alternative methods proposed in Section 2.4, when comparing the *PVs* of two and of three *BDTs* respectively, and both if the variance-covariance matrix is estimated through the delta method and if it is estimated through the bootstrap method. These simulation experiments were designed in a similar way to those carried out by Roldan-Nofuentes *et al.* [5], and consisted of the generation of 5000 random samples with multinomial distributions sized 50, 100, 200, 500, 1000, 2000 and 5000. For all of the study $\alpha = 5\%$ was set. All of the random samples were generated in such a way that in all of them it was possible to estimate the *PVs* and their variances-covariances. In the case of bootstrap, for each random sample 2000 samples with replacement were generated and from these $\hat{\tau}$ and $\hat{\Sigma}_{\hat{\tau}}$ were calculated. All of the random samples were generated from the *PVs* and the prevalence, without setting the values of sensitivity and specificity of each *BDT* in the following way:

1. As *PVs* we took the values $\{0.60, 0.65, \dots, 0.90, 0.95\}$, which are quite common values in clinical practice, and as values of the disease prevalence we took the values $\{0.05, 0.10, \dots, 0.90, 0.95\}$.
2. Once the *PVs* and the disease prevalence were set, the sensitivity and the specificity of each diagnostic test were calculated from equations (1.1), and then the maximum values of the dependence factors between the two *BDTs* were obtained from the values of the sensitivity and specificity of each diagnostic test applying the model of Vacek [20] for two *BDTs* and applying the model of Torrance-Rynard and Walter [21] for three *BDTs*. In Appendix A both models are summarized.
3. For two *BDTs*, as verification probabilities we took the values

$$(\lambda_{11} = 0.70, \lambda_{10} = \lambda_{01} = 0.40, \lambda_{00} = 0.10)$$

and

$$(\lambda_{11} = 0.95, \lambda_{10} = \lambda_{01} = 0.60, \lambda_{00} = 0.30),$$

which can be considered a scenario with low verification and a scenario

with high verification respectively. For three *BDTs*, we took the values

$$(\lambda_{111} = 0.70, \lambda_{110} = 0.40, \lambda_{101} = 0.40, \lambda_{100} = 0.25, \\ \lambda_{011} = 0.40, \lambda_{010} = 0.25, \lambda_{001} = 0.25, \lambda_{000} = 0.05)$$

and

$$(\lambda_{111} = 1, \lambda_{110} = 0.80, \lambda_{101} = 0.80, \lambda_{100} = 0.40, \\ \lambda_{011} = 0.80, \lambda_{010} = 0.40, \lambda_{001} = 0.40, \lambda_{000} = 0.20),$$

which can also be considered as scenarios with low and high verification.

4. In the case of two *BDTs*, the probabilities of the multinomial distributions were calculated from the equations of the model of Vacek [20] (Appendix A). In the case of three *BDTs*, the probabilities of the multinomial distributions were calculated from the model of Torrance-Rynard and Walter [21] (Appendix A).

3.1. Two *BDTs*

In Table 1 we show the results obtained for the type I errors probabilities and the powers when comparing the *PVs* of two *BDTs*, for different values of the *PVs* and for intermediate and high dependence factors, when the *PVs* are estimated through maximum likelihood and the variance-covariance matrix is estimated through the delta method (other tables with results from the simulation experiments can be requested from the authors). Regarding the type I error probability, the global hypothesis test has a type I error probability which, in general terms, fluctuates around the nominal error of 5% especially when $n \geq 500$. In some cases, especially when $n \leq 200$ and the verification probabilities are low and/or the dependence factors are high, the type I error probability may overwhelm the nominal error. This may be due to the fact that the samples are not large enough, and therefore some frequencies of the multinomial distribution which are equal to zero, and the variance-covariance matrix are not well represented. Regarding Methods 1, 2 and 3 (Bonferroni, Holm and Hochberg), the type I error probability of each one of them performs in a similar way to that of the global test, although it is usually somewhat lower than the nominal error (especially for $n \geq 2000$).

Regarding the power, in general it is necessary to have samples of between 500 and 1000 individuals (depending on the verification probabilities) so that the power of the global hypothesis test is high (higher than 80% or 90%). The power of the global hypothesis test increases when there is an increase in the verification probabilities; whereas the increase in the dependence factors does not have a clear effect on the power of the global hypothesis test (sometimes it increases and sometimes it decreases). Regarding Methods 1, 2 and 3, their respective powers perform in a similar way to that of the global test, although the power of each one of them is slightly lower than that of the global test.

Regarding the solution of the global test applying the bootstrap method, the results obtained are almost identical to those obtained through the method of maximum likelihood and the delta method. Therefore, in terms of the type I error probability and the power there is practically no difference between solving the global hypothesis test through the maximum likelihood method and the bootstrap method, although the bootstrap requires a greater computational effort.

3.2. Three *BDTs*

In Table 2 we show some of the results obtained for the type I error probability and the power when comparing the *PVs* of three *BDTs*, also for different values of the *PVs* and for intermediate and high dependence factors, when the *PVs* are estimated through maximum likelihood and the variance-covariance matrix is estimated through the delta method (other tables with results from the simulation experiments can be requested from the authors). For three *BDTs* we have not considered sample sizes smaller than 100, since with smaller samples there are too many frequencies equal to 0 (above all when the prevalence is low and/or the verification probabilities are low) and it is not possible to calculate the estimators or the variances-covariances. In general terms, the conclusions reached are similar to those obtained for two *BDTs*, although for the global test and for methods 1, 2 and 3 it is necessary to have larger sample sizes so that the type I error probability fluctuates around the nominal error.

With regard to the power of each method, this increases with an increase in the verification probabilities, and decreases when there is an increase in the values of the dependence factors. In very general terms, when the verification probabilities are low it is necessary to have samples of between 500 and 1000 individuals so that the power of the global test is higher than 80% or 90% (depending on the values of the dependence factors), although in some situations (high dependence factors) it is necessary to have very large samples ($n \geq 5000$) in order to reach this power. Regarding Methods 1, 2 and 3, in general terms there is no important difference in power in relation to the global hypothesis test, especially when $n \geq 500$, whilst for smaller sample sizes the global test is somewhat more powerful than for the other three methods.

3.3. Conclusions

From the analysis of the results obtained in the simulation experiments one may conclude that the global hypothesis test based on the chi-squared distribution displays the performance of an asymptotic hypothesis test (from a certain sample size onwards, its type I error probability fluctuates around the nominal error).

In general terms, its type I error probability fluctuates around the nominal error (especially for $n \geq 500$) and it is necessary to have large samples ($n \geq 500$) so that the power is greater than 80%. From the results obtained in the simulation experiments carried out, the global hypothesis test

$$H_0: PPV_1 = PPV_2 = \dots = PPV_J \quad \text{and} \quad NPV_1 = NPV_2 = \dots = NPV_J ,$$

$$H_1: \text{at least one equality is not true} ,$$

can be solved through the following procedure:

1. Solving the global hypothesis test based on the chi-squared distribution to a global error rate of α using the statistics given by equations (2.4) or bootstrap method.
2. If the global hypothesis test is not significant, then one cannot reject the homogeneity of the J $PPVs$ and of the J $NPVs$. If the global hypothesis test is significant to an error rate of α , in order to investigate the causes of the significance the following marginal hypothesis tests are solved

$$H_0: PPV_i = PPV_j \quad \text{vs} \quad H_1: PPV_i \neq PPV_j$$

and

$$H_0: NPV_i = NPV_j \quad \text{vs} \quad H_1: NPV_i \neq NPV_j$$

using the statistics given by equation (2.5), and applying some of the methods of multiple comparison used (Bonferroni, Holm or Hochberg) to an error rate of α .

4. EXAMPLE

The results obtained in Section 2 and the procedure given in Section 3.3 were applied to the diagnosis of coronary stenosis. Coronary stenosis is a disease that consists of the obstruction of the coronary artery and its diagnosis can be made through a dobutamine echocardiogram, a stress echocardiogram or a CT scan, and as the gold standard a coronary angiography is used. Coronary angiography may cause different reactions in patients (thrombosis, heart attacks, infections, even death) and therefore not all patients are verified with the gold standard. In Table 3 (Study of coronary stenosis), we show the results obtained by applying the dobutamine echocardiogram (variable T_1), the stress echocardiogram (variable T_2) and the CT scan (variable T_2) to a sample of 2455 males over 45 years of age and by only applying the coronary angiography (variable D) to a subset of these individuals. This study was carried out in two phases: firstly, the three $BDTs$ were applied to all of the individuals in the sample, and secondly the gold standard was applied to a subset of these individuals depending on only the results of the three diagnostic tests. This data are part of a study carried out at the University Hospital in Granada (Spain). In this example, one can assume that

the missing data mechanism is ignorable, and therefore the results from Section 3 can be applied. The values of the estimators of the PVs are $\widehat{PPV}_1 = 0.742$, $\widehat{PPV}_2 = 0.622$, $\widehat{PPV}_3 = 0.805$, $\widehat{NPV}_1 = 0.933$, $\widehat{NPV}_2 = 0.850$, $\widehat{NPV}_3 = 0.952$, and applying the delta method, the estimated asymptotic variance-covariance matrix is

$$\hat{\Sigma}_{\hat{\tau}} = \begin{pmatrix} 0.000234 & 0.000108 & 0.000086 & 0 & -0.000063 & -0.000038 \\ 0.000108 & 0.000258 & 0.000106 & -0.000035 & 0 & -0.000025 \\ 0.000086 & 0.000106 & 0.0000239 & -0.000059 & -0.000069 & 0 \\ 0 & -0.000034 & -0.000059 & 0.000114 & 0.000080 & 0.000045 \\ 0.000063 & 0 & 0.000069 & 0.000080 & 0.000169 & 0.000064 \\ 0.000038 & 0.000025 & 0 & 0.000045 & 0.000064 & 0.000085 \end{pmatrix}.$$

Applying equation (2.4) it holds that $Q^2 = 145.103$ (p -value = 0), and therefore we reject the equality of the three PPVs and of the three NPVs. In order to investigate the causes of the significance, the marginal hypothesis tests ($H_0: PPV_i = PPV_j$ and $H_0: NPV_i = NPV_j$) are solved. In Table 3 (Marginal hypothesis tests), we show the results obtained for each one of the six hypothesis tests that compare the PVs. Applying the Bonferroni method, the Holm method or the Hochberg method, it holds that the three PPVs are different, and that the PPV of the CT scan is the largest, followed by that of the dobutamine echocardiogram and, finally, that of the stress echocardiogram.

Table 3: Data from the study of coronary stenosis and marginal hypothesis tests.

Study of coronary stenosis									
	$T_1 = 1$				$T_1 = 0$				
	$T_2 = 1$		$T_2 = 0$		$T_2 = 1$		$T_2 = 0$		
	$T_3 = 1$	$T_3 = 0$	Total						
$V = 1$									
$D = 1$	457	30	84	5	34	0	7	1	618
$D = 0$	41	23	5	61	16	86	32	95	359
$V = 0$	92	31	85	120	42	195	88	825	1478
Total	590	84	174	186	92	281	127	921	2455

Marginal hypothesis tests		
Hypothesis test	z	Two sided p -value
$H_0: PPV_1 = PPV_2$ vs $H_1: PPV_1 \neq PPV_2$	3.61	0.003
$H_0: PPV_1 = PPV_3$ vs $H_1: PPV_1 \neq PPV_3$	7.20	6.06×10^{-13}
$H_0: PPV_2 = PPV_3$ vs $H_1: PPV_2 \neq PPV_3$	10.79	0
$H_0: NPV_1 = NPV_2$ vs $H_1: NPV_1 \neq NPV_2$	7.46	8.37×10^{-14}
$H_0: NPV_1 = NPV_3$ vs $H_1: NPV_1 \neq NPV_3$	1.76	0.078
$H_0: NPV_2 = NPV_3$ vs $H_1: NPV_2 \neq NPV_3$	8.99	0

Regarding the *NPVs*, no significant differences were found between the *NPVs* of the dobutamine echocardiogram and of the *CT* scan, whilst the *NPV* of the dobutamine echocardiogram is significantly lower than the *NPVs* of the other two *BDTs*.

5. DISCUSSION

Different studies have examined the problem of the comparison of the *PVs* of two or more *BDTs* when the diagnostic tests and the gold standard are applied to all of the individuals in a random sample. These models cannot be applied when a subset of individuals in the random sample have not had their disease status verified through the application of the gold standard, since the results obtained may be affected by the verification bias. In this article, we have studied a global hypothesis test to simultaneously compare the *PVs* of two or more *BDTs* when for a subset of individuals in the sample the disease status (either present or absent) is unknown. The global hypothesis test is based on the chi-squared distribution, and can be solved through the method of maximum likelihood and the delta method (equation (2.4) or through the bootstrap method, although the latter requires a greater computational effort. In terms of the type I error probability, both methods lead to very similar results, and the type I error probability fluctuates around the nominal error especially for $n \geq 500$. Other alternative methods to solve the global hypothesis test have been studied. The method based on the marginal comparisons of the *PPVs* (*NPVs*) to an error rate of $\alpha = 5\%$ leads to a type I error probability that clearly overwhelms the nominal error, and therefore this method may give rise to erroneous results. The methods based on marginal comparisons applying the corrections of Bonferroni, Holm and Hochberg respectively give rise to a type I error probability that fluctuates around the nominal error especially for $n \geq 500$. In terms of power, the global hypothesis test based on the chi-squared distribution (equation (2.4) or bootstrap method) is a little more powerful than the methods based on the corrections of Bonferroni, Holm and Hochberg respectively. Therefore, from the results of the simulation experiments carried out, the following method is proposed to compare the *PVs* of J *BDTs* in the presence of ignorable missing data: 1) Apply the global hypothesis test based on the chi-squared distribution to an error rate of α (equations (2.4) or bootstrap method); 2) If the global hypothesis test is significant to an error rate of α , investigate the causes of the significance solving the marginal hypothesis tests $H_0: PPV_i = PPV_j$ and $H_0: NPV_i = NPV_j$ along with a method of multiple comparisons (Bonferroni, Holm or Hochberg). This procedure is similar to that used in an analysis of variance. Firstly, the global test is solved and then a method of multiple comparisons is applied.

An alternative method to that proposed in Section 2 consists of solving the global test applying the Wilks method. Similar simulation experiments to those described in Section have demonstrated that the type I error probability and the power of this method are very similar to those obtained with the Wald method (equation (2.4)).

If all of the individuals are verified with the gold standard, and therefore all of the frequencies u_{i_1, \dots, i_J} are equal to 0, the method proposed by Roldán-Nofuentes *et al.* [5] is a particular case of the scenario analyzed in this study. Therefore, the simultaneous comparison of the PVs of two (or more) $BDTs$ in paired designs is a particular case of the scenario analyzed in this article.

APPENDIX A

In the case of two *BDTs*, the probabilities of the multinomial distribution were calculated applying the model of conditional dependence of Vacek [20], and their expressions are

$$\begin{aligned}\phi_{ij} &= \lambda_{ij} p \left\{ Se_1^i (1 - Se_1)^{1-i} Se_2^j (1 - Se_2)^{1-j} + \delta_{ij} \varepsilon_1 \right\}, \\ \varphi_{ij} &= \lambda_{ij} (1 - p) \left\{ Sp_1^{1-i} (1 - Sp_1)^i Sp_2^{1-j} (1 - Sp_2)^j + \delta_{ij} \varepsilon_0 \right\}, \\ \gamma_{ij} &= (1 - \lambda_{ij}) p \left\{ Se_1^i (1 - Se_1)^{1-i} Se_2^j (1 - Se_2)^{1-j} + \delta_{ij} \varepsilon_1 \right\} \\ &\quad + (1 - \lambda_{ij}) (1 - p) \left\{ Sp_1^{1-i} (1 - Sp_1)^i Sp_2^{1-j} (1 - Sp_2)^j + \delta_{ij} \varepsilon_0 \right\},\end{aligned}$$

where $\delta_{ij} = 1$ when $i = j$ and $\delta_{ij} = -1$, and ε_i is the dependence factor (covariance) between the two *BDTs* when $D = i$. In clinical practice, the two *BDTs* are usually conditionally dependent on the disease, and it is verified [20] that $0 < \varepsilon_1 < Se_1 (1 - Se_2)$ when $Se_2 > Se_1$ and $0 < \varepsilon_1 < Se_2 (1 - Se_1)$ when $Se_1 > Se_2$, and in the same way, $0 < \varepsilon_0 < Sp_1 (1 - Sp_2)$ when $Sp_2 > Sp_1$ and $0 < \varepsilon_0 < Sp_2 (1 - Sp_1)$ when $Sp_1 > Sp_2$. If the two *BDTs* are conditionally independent on the disease then $\varepsilon_1 = \varepsilon_0 = 0$.

In the case of three *BDTs*, the probabilities of the multinomial distributions were calculated applying the model of Torrance-Rynard and Walter [21]:

$$\begin{aligned}P(V = 1, D = 1, T_1 = i_1, T_2 = i_2, T_3 = i_3) &= \\ &= p \lambda_{i_1 i_2 i_3} \left\{ \prod_{j=1}^3 Se_j^{i_j} (1 - Se_j)^{1-i_j} + \sum_{j,k,j < k}^3 (-1)^{|i_j - i_k|} \delta_{jk} \right\},\end{aligned}$$

$$\begin{aligned}P(V = 1, D = 0, T_1 = i_1, T_2 = i_2, T_3 = i_3) &= \\ &= (1 - p) \lambda_{i_1 i_2 i_3} \left\{ \prod_{j=1}^3 Sp_j^{1-i_j} (1 - Sp_j)^{i_j} + \sum_{j,k,j < k}^3 (-1)^{|i_j - i_k|} \varepsilon_{jk} \right\}\end{aligned}$$

and

$$\begin{aligned}P(V = 0, T_1 = i_1, T_2 = i_2, T_3 = i_3) &= \\ &= p (1 - \lambda_{i_1 i_2 i_3}) \left\{ \prod_{j=1}^3 Se_j^{i_j} (1 - Se_j)^{1-i_j} + \sum_{j,k,j < k}^3 (-1)^{|i_j - i_k|} \delta_{jk} \right\} \\ &\quad + (1 - p) (1 - \lambda_{i_1 i_2 i_3}) \left\{ \prod_{j=1}^3 Sp_j^{1-i_j} (1 - Sp_j)^{i_j} + \sum_{j,k,j < k}^3 (-1)^{|i_j - i_k|} \varepsilon_{jk} \right\},\end{aligned}$$

with $i_j = 0, 1$, $i_k = 0, 1$ and $j, k = 1, 2, 3$, and where δ_{jk} (ε_{jk}) is the factor of dependence between the j -th *BDT* and k -th *BDT* when $D = 1$ ($D = 0$).

The factors of dependence δ_{jk} and ε_{jk} verify restrictions that depend on the values of sensitivity and specificity of the three *BDTs*. In order to simplify the simulation experiments, it has been considered that $\delta_{ij} = \delta$ and $\varepsilon_{ij} = \varepsilon$, so that the factors of dependence verify the following restrictions:

$$\delta \leq \text{Min}\left\{(1-Se_1)(1-Se_2)Se_3, (1-Se_1)Se_2(1-Se_3), Se_1(1-Se_2)(1-Se_3)\right\}$$

and

$$\varepsilon \leq \text{Min}\left\{(1-Sp_1)(1-Sp_2)Sp_3, (1-Sp_1)Sp_2(1-Sp_3), Sp_1(1-Sp_2)(1-Sp_3)\right\}.$$

In clinical practice, factors δ_{jk} and/or ε_{jk} are greater than zero, so that the *BDTs* are conditionally dependent on the disease status. When $\delta_{jk} = \varepsilon_{jk} = 0$ the three *BDTs* are conditionally independent on the disease status.

ACKNOWLEDGMENTS

This research was supported by the Spanish Ministry of Economy, Grant Number MTM2012-35591. We thank the two referees, the Associate Editor and the Editor (M. Ivette Gomes) of *REVSTAT* for their helpful comments that improved the quality of the paper.

REFERENCES

- [1] BENNET, B. M. (1972). On comparison of sensitivity, specificity and predictive value of a number of diagnostic procedures, *Biometrics*, **28**, 793–800.
- [2] BENNET, B. M. (1985). On tests for equality of predictive values for t diagnostic procedures, *Statistics in Medicine*, **4**, 525–539.
- [3] LEISENRING, W.; ALONZO, T. and PEPE, M. S. (2000). Comparisons of predictive values of binary medical diagnostic tests for paired designs, *Biometrics*, **56**, 354–351.
- [4] WNAG, W.; DAVIS, C. S. and SOONG, S. J. (2006). Comparison of predictive values of two diagnostic tests from the same sample of subjects using weighted least squares, *Statistics in Medicine*, **25**, 2215–2229.
- [5] ROLDÁN NOFUENTES, J. A.; LUNA DEL CASTILLO, J. D. and MONTERO ALONSO, M. A. (2012). Global hypothesis test to simultaneously compare the predictive values of two binary diagnostic tests, *Computational Statistics and Data Analysis, Special issue “Computational Statistics for Clinical Research”*, **56**, 1161–1173.

- [6] POLETO, F. Z.; SINGER, J. M. and PAULINO, C. D. (2011). Comparing diagnostic tests with missing data, *Journal of Applied Statistics*, **38**, 1207–1222.
- [7] BEGG, C. B. and GREENES, R. A. (1983). Assessment of diagnostic tests when disease verification is subject to selection bias, *Biometrics*, **39**, 207–215.
- [8] ZHOU, X. H. (1993). Maximum likelihood estimators of sensitivity and specificity corrected for verification bias, *Communication in Statistics Theory and Methods*, **22**, 3177–3198.
- [9] ROLDÁN NOFUENTES, J. A. and LUNA DEL CASTILLO, J. D. (2007). The effect of verification bias in the naive estimators of accuracy of a binary diagnostic test, *Communications in Statistics — Simulation and Computation*, **36**, 959–972.
- [10] ROLDÁN NOFUENTES, J. A. and LUNA DEL CASTILLO, J. D. (2008a). EM algorithm for comparing two binary diagnostic tests when not all the patients are verified, *Journal of Statistical Computation and Simulation*, **78**, 19–35.
- [11] ROLDÁN NOFUENTES, J. A. and LUNA DEL CASTILLO, J. D. (2008b). The effect of verification bias on the comparison of predictive values of two binary diagnostic tests, *Journal of Statistical Planning and Inference*, **138**, 959–963.
- [12] RUBIN, D. B. (1976). Inference and missing data, *Biometrika*, **4**, 73–89.
- [13] SCHAFER, J. L. (1997). *Analysis of Incomplete Multivariate Data*, Chapman and Hall/CRC, USA.
- [14] BISHOP, Y. M.; FIENBERG, S. E. and HOLLAND, P. W. (1975). *Discrete Multivariate Analysis: Theory and Practice*, The MIT Press.
- [15] AGRESTI, A. (2002). *Categorical Data Analysis*, John Wiley and Sons, New York.
- [16] EFRON, B. and TIBSHIRNI, R. J. (1993). *An Introduction to the Bootstrap*, Chapman and Hall, New York.
- [17] BONFERRONI, C. E. (1936). Teoria statistica delle classi e calcolo delle probabilità, *Pubblicazioni del R Istituto Superiore di Scienze Economiche e Commerciali di Firenze*, **8**, 3–62.
- [18] HOLM, S. (1979). A simple sequential rejective multiple testing procedure, *Scandinavian Journal of Statistics*, **6**, 65–70.
- [19] HOCHBERG, Y. (1988). A sharper Bonferroni procedure for multiple tests of significance, *Biometrika*, **75**, 800–802.
- [20] VACEK, P. (1985). The effect of conditional dependence on the evaluation of diagnostic tests, *Biometrics*, **23**, 959–968.
- [21] TORRANCE-RYNARD, V. L. and WALTER, S. D. (1997). Effects of dependent errors in the assessment of diagnostic test performance, *Statistics in Medicine*, **16**, 2157–2175.