

# Combdt: an R program to compare two binary diagnostic tests subject to a paired design

Jose Antonio Roldán-Nofuentes (✉ [jaroldan@ugr.es](mailto:jaroldan@ugr.es))

---

## Software

**Keywords:** binary diagnostic test; likelihood ratios; paired design; predictive values; Sensitivity and specificity.

**Posted Date:** April 1st, 2020

**DOI:** <https://doi.org/10.21203/rs.2.22525/v2>

**License:**  This work is licensed under a Creative Commons Attribution 4.0 International License.  
[Read Full License](#)

---

**Version of Record:** A version of this preprint was published at BMC Medical Research Methodology on June 5th, 2020. See the published version at <https://doi.org/10.1186/s12874-020-00988-y>.

1      **Compbdт: an R program to compare two binary diagnostic tests**

2                   **subject to a paired design**

4                   J. A. Roldán-Nofuentes

5                   Department of Statistics (Biostatistics), School of Medicine, University of Granada,

6                   Avenida de la Investigación 11, Granada, 18016, Spain

7                   Email: jaroldan@ugr.es

8      **Abstract**

9      **Background:** The comparison of the performance of two binary diagnostic tests is an  
10     important topic in Clinical Medicine. The most frequent type of sample design to compare  
11     two binary diagnostic tests is the paired design. This design consists of applying the two  
12     binary diagnostic tests to all of the individuals in a random sample, where the disease status of  
13     each individual is known through the application of a gold standard. This article presents an R  
14     program to compare parameters of two binary tests subject to a paired design.

15     **Results:** The “compbdт” program estimates the sensitivity and the specificity, the likelihood  
16     ratios and the predictive values of each diagnostic test applying the confidence intervals with  
17     the best asymptotic performance. The program compares the sensitivities and specificities of  
18     the two diagnostic tests simultaneously, as well as the likelihood ratios and the predictive  
19     values, applying the global hypothesis tests with the best performance in terms of type I error  
20     and power. When the global hypothesis test is significant, the causes of the significance are  
21     investigated solving the individual hypothesis tests and applying the multiple comparison  
22     method of Holm. The most optimal confidence intervals are also calculated for the difference  
23     or ratio between the respective parameters. Based on the data observed in the sample, the  
24     program also estimates the probability of making a type II error if the null hypothesis is not  
25     rejected, or estimates the power if the alternative hypothesis is accepted. The “compbdт”

26 program provides all the necessary results so that the researcher can easily interpret them. The  
27 estimation of the probability of making a type II error allows the researcher to decide about  
28 the reliability of the null hypothesis when this hypothesis is not rejected. The “compbdt”  
29 program has been applied to a real example on the diagnosis of coronary artery disease.

30 **Conclusions:** The “compbdt” program is one which is easy to use and allows the researcher  
31 to compare the most important parameters of two binary tests subject to a paired design. The  
32 “compbdt” program is available as supplementary material.

33 **Keywords:** binary diagnostic test; likelihood ratios; paired design; predictive values;  
34 Sensitivity and specificity.

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51    **Background**

52    A diagnostic test is a medical test that is applied to an individual in order to determine the  
53    presence or absence of a disease. When the result of a diagnostic test is positive or negative,  
54    the diagnostic test is called a binary diagnostic test. A stress test for the diagnosis of coronary  
55    disease is an example of binary diagnostic test. The performance of a binary diagnostic test is  
56    measured in terms of two fundamental parameters: sensitivity and specificity. The sensitivity  
57    ( $Se$ ) is the probability of the diagnostic test being positive when the individual has the disease,  
58    and the specificity ( $Sp$ ) is the probability of the diagnostic test being negative when the  
59    individual does not have it. The  $Se$  and the  $Sp$  of a diagnostic test are estimated in relation to a  
60    gold standard, which is a medical test which objectively determines whether or not an  
61    individual has the disease or not. An angiography for coronary disease is an example of a gold  
62    standard. Other parameters that are used to assess the performance of a diagnostic test are the  
63    likelihood ratios ( $LRs$ ) and the predictive values ( $PVs$ ) [1, 2]. When the diagnostic test is  
64    positive, the likelihood ratio, called the positive likelihood ratio ( $PLR$ ), is the ratio between  
65    the probability of correctly classifying an individual with the disease and the probability of  
66    incorrectly classifying an individual who does not have it, i.e.  $PLR = Se/(1 - Sp)$ . When the  
67    diagnostic test is negative, the likelihood ratio, called the negative likelihood ratio ( $NLR$ ), is  
68    the ratio between the probability of incorrectly classifying an individual who has the disease  
69    and the probability of correctly classifying an individual who does not have it, i.e.  
70     $NLR = (1 - Se)/Sp$ . The  $LRs$  only depend on  $Se$  and  $Sp$  of the diagnostic test and they are  
71    equivalent to a relative risk. The positive predictive value ( $PPV$ ) is the probability of an  
72    individual having the disease when the result of the diagnostic test is positive, and the  
73    negative predictive value ( $NPV$ ) is the probability of an individual not having the disease  
74    when the result of the diagnostic test is negative. The  $PVs$  represent the accuracy of the  
75    diagnostic test when it is applied to a cohort of individuals, and they are measures of the

76 clinical accuracy of the diagnostic test. The *PVs* depend on the *Se* and the *Sp* of the diagnostic  
77 test and on the disease prevalence (*p*), and are easily calculated applying Bayes' Theorem i.e.

78 
$$PPV = \frac{p \times Se}{p \times Se + (1-p) \times (1-Sp)} \text{ and } NPV = \frac{(1-p) \times Sp}{p \times (1-Se) + (1-p) \times Sp}.$$

79 Whereas the *Se* and the *Sp* quantify how well the diagnostic test reflects the true disease status  
80 (present or absent), the *PVs* quantify the clinical value of the diagnostic test, since both the  
81 individual and the clinician are more interested in knowing how probable it is to have the  
82 disease given a diagnostic test result.

83 The comparison of the performance of two diagnostic tests with respect to a gold standard  
84 is an important topic in Clinical Medicine and Epidemiology. The most frequent type of  
85 sample design to compare two diagnostic tests with respect to a gold standard is paired design  
86 [1, 2]. This design consists of applying the two diagnostic tests, *Test 1* and *Test 2*, to all of the  
87 individuals in a random sample sized *n*, where the disease status of each individual is known  
88 through the application of a gold standard. Therefore, subject to a paired design the two  
89 diagnostic tests and the gold standard are applied to all of the individuals in a single random  
90 sample, whose size (*n*) has been set by the researcher. Paired design is the most efficient type  
91 of design to compare two binary diagnostic tests as it minimizes the impact of the between-  
92 individual variability, therefore this manuscript focuses on paired design. The comparison of  
93 two diagnostic tests subject to this type of design leads to the frequencies that are shown in  
94 Table 1, where  $s_{ij}$  ( $r_{ij}$ ) be the number of diseased (non-diseased) patients in which the *Test 1*  
95 gives a result *i* (1 positive and 0 negative) and *Test 2* gives a result *j* (1 positive and 0  
96 negative).

97

98

99

100

Table 1. Frequencies subject to a paired design.

	Test 1 positive		Test 1 negative		Total
	Test 2 positive	Test 2 negative	Test 2 positive	Test 2 negative	
Disease	$s_{11}$	$s_{10}$	$s_{01}$	$s_{00}$	$s$
No disease	$r_{11}$	$r_{10}$	$r_{01}$	$r_{00}$	$r$
Total	$n_{11}$	$n_{10}$	$n_{01}$	$n_{00}$	$n$

101

102 This article presents a program called “compbdt” (Comparison of two Binary Diagnostic  
 103 Tests) written in R [3] which allows us to estimate and compare the performance (measured  
 104 in terms of the previous parameters) of two diagnostic tests subject to a paired design  
 105 applying the statistical methods with the best asymptotic performance, i.e. for the confidence  
 106 intervals we used the intervals that have a better coverage and average width, and for the  
 107 hypothesis tests we used the methods that have the best behaviour in terms of type I error and  
 108 power. In the next section, the methods of estimation and of comparison of the parameters are  
 109 summarized, and the “compbdt” program is explained. The results are applied to a real  
 110 example of the diagnosis of coronary artery disease, and finally some conclusions are given.

111

## 112 **Implementation**

113 The estimation and comparison of parameters of two diagnostic tests has been the subject of  
 114 numerous studies in Statistics literature. We will now describe the statistical methods  
 115 implemented in the “compbdt” program to estimate the parameters and to compare the  
 116 respective parameters subject to a paired design. The methods used are those that have a  
 117 better asymptotic behaviour in terms of coverage for the confidence intervals and in terms of  
 118 type I error and power for hypothesis tests.

119

## 120 **Estimation of the parameters**

121 The estimation of the sensitivity, the specificity and the predictive values of each diagnostic  
 122 test consists of the estimation of a binomial proportion. There are numerous confidence

123 intervals proposed to estimate a binomial proportion. Yu et al [4] proposed a new interval,  
 124 based on a modification of the Wilson interval, to estimate a binomial proportion,  
 125 demonstrating that this interval shows a better asymptotic performance than the rest of the  
 126 existing intervals. For the sensitivity of each diagnostic test, the estimators are

$$127 \quad \hat{Se}_1 = \frac{s_{11} + s_{10}}{s} \text{ and } \hat{Se}_2 = \frac{s_{11} + s_{01}}{s},$$

128 and their standard errors ( $SE$ ) are

$$129 \quad SE(\hat{Se}_i) = \sqrt{\frac{\hat{Se}_i(1-\hat{Se}_i)}{n\hat{p}}},$$

130 with  $i=1,2$ , and where  $\hat{p}=s/n$  is the estimator of the disease prevalence. The Yu et al  
 131 confidence interval for sensitivity  $Se_i$ , with  $i=1,2$ , is

$$132 \quad Se_i \in 0.5 + \frac{s + z_{1-\alpha/2}^4 / 53}{s + z_{1-\alpha/2}^4} (\hat{Se}_i - 0.5) \pm \frac{z_{1-\alpha/2}}{s + z_{1-\alpha/2}^2} \sqrt{s(1-\hat{Se}_i)\hat{Se}_i + \frac{z_{1-\alpha/2}^2}{4}},$$

133 where  $z_{1-\alpha/2}$  is the  $100(1-\alpha/2)th$  percentile of the standard normal distribution. For the  
 134 specificities, the estimators are

$$135 \quad \hat{Sp}_1 = \frac{r_{01} + r_{00}}{r} \text{ and } \hat{Sp}_2 = \frac{r_{10} + r_{00}}{r},$$

136 and their standard errors ( $SE$ ) are

$$137 \quad SE(\hat{Sp}_i) = \sqrt{\frac{\hat{Sp}_i(1-\hat{Sp}_i)}{n(1-\hat{p})}}$$

138 The intervals for the specificities are obtained analogously by replacing  $\hat{Se}_i$  with  $\hat{Sp}_i$  and  $s$   
 139 with  $r$ .

140 For the predictive values, the estimators of the  $PPVs$  are

$$141 \quad \hat{PPV}_1 = \frac{s_{10} + s_{11}}{s_{10} + s_{11} + r_{10} + r_{11}} \text{ and } \hat{PPV}_2 = \frac{s_{01} + s_{11}}{s_{01} + s_{11} + r_{01} + r_{11}},$$

142 and their standard errors are

$$143 \quad SE(\hat{PPV}_1) = \sqrt{\frac{(s_{10} + s_{11})(r_{10} + r_{11})}{n(s_{10} + s_{11} + r_{10} + r_{11})^3}} \text{ and } SE(\hat{PPV}_2) = \sqrt{\frac{(s_{01} + s_{11})(r_{01} + r_{11})}{n(s_{01} + s_{11} + r_{01} + r_{11})^3}}.$$

144 The estimators of the *NPVs* are

$$145 \quad \hat{NPV}_1 = \frac{r_{00} + r_{01}}{s_{00} + s_{01} + r_{00} + r_{01}} \text{ and } \hat{NPV}_2 = \frac{r_{00} + r_{10}}{s_{00} + s_{10} + r_{00} + r_{10}},$$

146 and their standard errors are

$$147 \quad SE(\hat{NPV}_1) = \sqrt{\frac{(s_{00} + s_{01})(r_{00} + r_{01})}{n(s_{00} + s_{01} + r_{00} + r_{01})^3}} \text{ and } SE(\hat{NPV}_2) = \sqrt{\frac{(s_{00} + s_{10})(r_{00} + r_{10})}{n(s_{00} + s_{10} + r_{00} + r_{10})^3}}.$$

148 For *PPV* and *NPV* of *Test 1*, the Yu et al confidence intervals are

$$149 \quad 0.5 + \frac{n_{1\Box} + z_{1-\alpha/2}^4 / 53}{n_{1\Box} + z_{1-\alpha/2}^4} (\hat{PPV}_1 - 0.5) \pm \frac{z_{1-\alpha/2}}{n_{1\Box} + z_{1-\alpha/2}^2} \sqrt{n_{1\Box} (1 - \hat{PPV}_1) \hat{PPV}_1 + \frac{z_{1-\alpha/2}^2}{4}}$$

150 and

$$151 \quad 0.5 + \frac{n_{0\Box} + z_{1-\alpha/2}^4 / 53}{n_{0\Box} + z_{1-\alpha/2}^4} (\hat{NPV}_1 - 0.5) \pm \frac{z_{1-\alpha/2}}{n_{0\Box} + z_{1-\alpha/2}^2} \sqrt{n_{0\Box} (1 - \hat{NPV}_1) \hat{NPV}_1 + \frac{z_{1-\alpha/2}^2}{4}}$$

152 where  $n_{1\Box} = (n_{11} + n_{10})$  and  $n_{0\Box} = n_{01} + n_{00}$ , respectively. The confidence intervals for *PPV* and

153 *NPV* of *Test 2* are obtained analogously by replacing  $n_{1\Box}$  with  $n_{1\Box} = n_{11} + n_{01}$  and  $\hat{PPV}_1$  with

154  $\hat{PPV}_2$ , and replacing  $n_{0\Box}$  with  $n_{0\Box} = n_{10} + n_{00}$  and  $\hat{NPV}_1$  with  $\hat{NPV}_2$ , respectively.

155 Regarding the likelihood ratios, the estimators of *PLRs* are

$$156 \quad \hat{PLR}_1 = \frac{r(s_{11} + s_{10})}{s(r_{11} + r_{10})} \text{ and } \hat{PLR}_2 = \frac{r(s_{11} + s_{01})}{s(r_{11} + r_{01})},$$

157 and their standard errors are

$$158 \quad SE(\hat{PLR}_i) = \sqrt{\frac{\hat{Se}_i^2 \times \hat{Var}(\hat{Sp}_i) + (1 - \hat{Sp}_i)^2 \times \hat{Var}(\hat{Se}_i)}{(1 - \hat{Sp}_i)^4}}, \quad i = 1, 2,$$

159 where  $\hat{Var}(\hat{Se}_i) = [SE(\hat{Se}_i)]^2$  and  $\hat{Var}(\hat{Sp}_i) = [SE(\hat{Sp}_i)]^2$ . The estimators of *NLRs* are

160 
$$\hat{NLR}_1 = \frac{r(s_{01} + s_{00})}{s(r_{01} + r_{00})} \text{ and } \hat{NLR}_2 = \frac{r(s_{10} + s_{00})}{s(r_{10} + r_{00})},$$

161 and their standard errors are

162 
$$SE(\hat{NLR}_i) = \sqrt{\frac{(1 - \hat{Se}_i)^2 \times \hat{Var}(\hat{Sp}_i) + \hat{Sp}_i^2 \times \hat{Var}(\hat{Se}_i)}{\hat{Sp}_i^4}}, \quad i = 1, 2.$$

163 The *LRs* are the ratio of two independent binomial proportions, i.e. a relative risk. Martín-  
 164 Andrés and Álvarez-Hernández [5] compared seventy-three confidence intervals for the ratio  
 165 of two independent binomial proportions, and concluded that the interval with the best  
 166 performance is the interval based on an approximation to the score method adding 0.5 to the  
 167 observed frequencies. For *Test 1*, these confidence intervals are:

168 
$$PLR_1 \in \frac{\tilde{n}\tilde{s}_{1\Box}\tilde{r}_{1\Box} + \frac{z_{1-\alpha/2}^2}{2}(\tilde{s}\tilde{s}_{1\Box} + \tilde{r}\tilde{r}_{1\Box}' - 2\tilde{s}_{1\Box}\tilde{r}_{1\Box}) \pm z_{1-\alpha/2}\sqrt{\frac{\tilde{n}^2\tilde{s}_{1\Box}\tilde{r}_{1\Box}}{\tilde{r}_{1\Box}}[\tilde{s}_{1\Box} + \tilde{r}_{1\Box} - \tilde{n}\tilde{Se}_1(1 - \tilde{Sp}_1)] + \frac{z_{1-\alpha/2}^2}{4}(\tilde{s}\tilde{s}_{1\Box} - \tilde{r}\tilde{r}_{1\Box})^2}}{\tilde{r}_{1\Box}[\tilde{n}\tilde{s}(1 - \tilde{Sp}_1) - z_{1-\alpha/2}^2(\tilde{s} - \tilde{r}_{1\Box})]}$$

169 and

170 
$$NLR_1 \in \frac{\tilde{n}\tilde{s}_{0\Box}\tilde{r}_{0\Box} + \frac{z_{1-\alpha/2}^2}{2}(\tilde{s}\tilde{s}_{0\Box} + \tilde{r}\tilde{r}_{0\Box}' - 2\tilde{s}_{0\Box}\tilde{r}_{0\Box}) \pm z_{1-\alpha/2}\sqrt{\frac{\tilde{n}^2\tilde{s}_{0\Box}\tilde{r}_{0\Box}}{\tilde{r}_{0\Box}}[\tilde{s}_{0\Box} + \tilde{r}_{0\Box} - \tilde{n}(1 - \tilde{Se}_1)\tilde{Sp}_1] + \frac{z_{1-\alpha/2}^2}{4}(\tilde{s}\tilde{s}_{0\Box} - \tilde{r}\tilde{r}_{0\Box})^2}}{\tilde{r}_{0\Box}[\tilde{n}\tilde{s}\tilde{Sp}_1 - z_{1-\alpha/2}^2(\tilde{s} - \tilde{r}_{0\Box})]}$$

171 where  $\tilde{s}_{1\Box} = s_{1\Box} + 0.5$ ,  $\tilde{s}_{0\Box} = s_{0\Box} + 0.5$ ,  $\tilde{r}_{1\Box} = r_{1\Box} + 0.5$ ,  $\tilde{r}_{0\Box} = r_{0\Box} + 0.5$ ,  $\tilde{s} = s + 1$ ,  $\tilde{r} = r + 1$ ,  $\tilde{n} = n + 2$ ,

172  $\tilde{Se}_1 = \tilde{s}_{1\Box}/\tilde{s}$  and  $\tilde{Sp}_1 = \tilde{r}_{1\Box}/\tilde{r}$ . If the lower limit of the interval for  $PLR_1$  is less than  $\tilde{s}_{1\Box}/(\tilde{n} - \tilde{r}_{1\Box})$

173 or greater than  $\hat{PLR}_1$ , then the lower limit of the confidence interval is

174 
$$\frac{\tilde{s}_{1\Box}(1 - \tilde{Sp}_1) + \frac{z_{1-\alpha/2}^2}{2} - z_{1-\alpha/2}\sqrt{\frac{z_{1-\alpha/2}^2}{4} + \tilde{s}_{1\Box}(1 - \tilde{Sp}_1 - \tilde{Se}_2)}}{\tilde{s}(1 - \tilde{Sp}_1)^2 + z_{1-\alpha/2}^2},$$

175 and if the upper limit of this interval is greater than  $(\tilde{n} - \tilde{s}_{1\square})/\tilde{r}_{1\square}$  or lower than  $\hat{PLR}_1$ , then the  
 176 upper limit of the confidence interval is

$$177 \quad \frac{\tilde{r}_{1\square}\tilde{S}e_1 + \frac{z_{1-\alpha/2}^2}{2} + z_{1-\alpha/2}\sqrt{\frac{z_{1-\alpha/2}^2}{4} + \tilde{r}_{1\square}(\tilde{S}e_1 + \tilde{S}p_1 - 1)}}{\tilde{r}(1 - \tilde{S}p_1)^2}.$$

178 Regarding the confidence interval for  $NLR_1$ , if the lower limit of this interval is less than  
 179  $\tilde{s}_{0\square}/(\tilde{n} - \tilde{r}_{0\square})$  or greater than  $\hat{NLR}_1$ , then the lower limit of the confidence interval is

$$180 \quad \frac{\tilde{s}_{0\square}\tilde{S}p_1 + \frac{z_{1-\alpha/2}^2}{2} - z_{1-\alpha/2}\sqrt{\frac{z_{1-\alpha/2}^2}{4} + \tilde{s}_{0\square}(\tilde{S}p_1 + \tilde{S}e_1 - 1)}}{\tilde{s}\tilde{S}p_1^2 + z_{1-\alpha/2}^2},$$

181 and if the upper limit of this interval is greater than  $(\tilde{n} - \tilde{s}_{0\square})/\tilde{r}_{0\square}$  or less than  $\hat{NLR}_1$ , then the  
 182 upper limit of the confidence interval is

$$183 \quad \frac{\tilde{r}_{0\square}(1 - \tilde{S}e_1) + \frac{z_{1-\alpha/2}^2}{2} + z_{1-\alpha/2}\sqrt{\frac{z_{1-\alpha/2}^2}{4} + \tilde{r}_{0\square}(1 - \tilde{S}e_1 - \tilde{S}p_1)}}{\tilde{r}\tilde{S}p_1^2}.$$

184 The confidence intervals for *L*<sub>2</sub>s of *Test 2* are obtained analogously by replacing  $\tilde{s}_{1\square}$  with  
 185  $\tilde{s}_{1\square} = s_{1\square} + 0.5$ ,  $\tilde{r}_{1\square}$  with  $\tilde{r}_{1\square} = r_{1\square} + 0.5$ ,  $\tilde{s}_{0\square}$  with  $\tilde{s}_{0\square} = s_{0\square} + 0.5$ ,  $\tilde{r}_{0\square}$  with  $\tilde{r}_{0\square} = r_{0\square} + 0.5$ ,  $\tilde{S}e_1$  with  
 186  $\tilde{S}e_2 = \tilde{s}_{1\square}/\tilde{s}$  and  $\tilde{S}p_1$  with  $\tilde{S}p_2 = \tilde{r}_{1\square}/\tilde{r}$ .

187 The “compbdt” program also estimates the prevalence of the disease. The estimator of the  
 188 prevalence is  $\hat{p} = s/n$ , the standard error is  $\sqrt{\hat{p}(1 - \hat{p})/n}$  and the Yu et al confidence interval  
 189 for the prevalence is

$$190 \quad p \in 0.5 + \frac{n + z_{1-\alpha/2}^4/53}{n + z_{1-\alpha/2}^4}(\hat{p} - 0.5) \pm \frac{z_{1-\alpha/2}}{n + z_{1-\alpha/2}^2}\sqrt{n(1 - \hat{p})\hat{p} + \frac{z_{1-\alpha/2}^2}{4}}.$$

191

192

193    **Comparison of the parameters**

194    The comparison of parameters of two diagnostic tests subject to a paired design has been the  
195    subject of different studies. The hypothesis tests with the best performance, in terms of type I  
196    and power error, to compare the parameters of two diagnostic tests are presented below.

197

198    *Comparison of the sensitivities and the specificities*

199    Traditionally, the comparison of two sensitivities and of two specificities was carried out  
200    solving the hypothesis tests  $H_0 : Se_1 = Se_2$  vs  $H_1 : Se_1 \neq Se_2$  and  $H_0 : Sp_1 = Sp_2$  vs  
201     $H_1 : Sp_1 \neq Sp_2$  each one of them to an  $\alpha$  error, applying a comparison test of two paired  
202    binomial proportions (e.g. the McNemar test) [2]. Recently, Roldán-Nofuentes and Sidaty-  
203    Regad [6] have studied different methods to compare the two sensitivities and the two  
204    specificities individually and also simultaneously, and carried out simulation experiments to  
205    compare these methods. The results of the simulation experiments showed that disease  
206    prevalence and sample size have an important effect on the type I errors and powers of the  
207    methods analysed, and from the results obtained some general rules of application were given  
208    in terms of the prevalence and the sample size. These rules are:

209    a). When the prevalence is small ( $\leq 10\%$ ) and the sample size  $n$  is  $\leq 100$ , solve the tests

210     $H_0 : Se_1 = Se_2$  and  $H_0 : Sp_1 = Sp_2$  individually applying the Wald test (or the likelihood ratio  
211    test) along with the Bonferroni or Holm method to an  $\alpha$  error. However, the second method  
212    has the disadvantage that it can only be applied if the frequencies of the discordant pairs are  
213    greater than zero. For  $H_0 : Se_1 = Se_2$  the Wald test statistic is

214    
$$\chi_{wSe}^2 = \frac{s(s_{10} - s_{01})^2}{4s_{10}s_{01} + (s_{11} + s_{00})(s_{10} + s_{01})},$$

215    and for  $H_0 : Sp_1 = Sp_2$  the Wald test statistic is

216

$$\chi_{WSp}^2 = \frac{r(r_{10} - r_{01})^2}{4r_{10}r_{01} + (r_{11} + r_{00})(r_{10} + r_{01})}.$$

217 Likelihood ratio test statistics are

218

$$\chi_{LRTSe}^2 = 2 \left[ s_{10} \ln \left( \frac{2s_{10}}{s_{10} + s_{01}} \right) + s_{01} \ln \left( \frac{2s_{01}}{s_{10} + s_{01}} \right) \right]$$

219 and

220

$$\chi_{LRSP}^2 = 2 \left[ r_{10} \ln \left( \frac{2r_{10}}{r_{10} + r_{01}} \right) + r_{01} \ln \left( \frac{2r_{01}}{r_{10} + r_{01}} \right) \right],$$

221 respectively. These statistics have a standard normal distribution. Both methods, the Wald test  
 222 and the likelihood ratio test, have a very similar asymptotic performance. However, the  
 223 second method has the disadvantage that it can only be applied if the frequencies of the  
 224 discordant pairs are greater than zero.

225 b) In any other situation, solve the global test  $H_0 : (Se_1 = Se_2 \text{ and } Sp_1 = Sp_2)$  vs  
 226  $H_1 : (Se_1 \neq Se_2 \text{ and/or } Sp_1 \neq Sp_2)$  to an  $\alpha$  error applying the Wald test or the likelihood ratio  
 227 test, i.e.

228

$$\chi_w^2 = \frac{s(s_{10} - s_{01})^2}{4s_{10}s_{01} + (s_{11} + s_{00})(s_{10} + s_{01})} + \frac{r(r_{10} - r_{01})^2}{4r_{10}r_{01} + (r_{11} + r_{00})(r_{10} + r_{01})}$$

229 and

230

$$\chi_{LRT}^2 = 2 \left[ s_{10} \ln \left( \frac{2s_{10}}{s_{10} + s_{01}} \right) + s_{01} \ln \left( \frac{2s_{01}}{s_{10} + s_{01}} \right) + r_{10} \ln \left( \frac{2r_{10}}{r_{10} + r_{01}} \right) + r_{01} \ln \left( \frac{2r_{01}}{r_{10} + r_{01}} \right) \right].$$

231 The distribution of both statistics is a chi-square with two degrees of freedom when the null  
 232 hypothesis is true. In this situation, if the global test is not significant then the equality of the  
 233 accuracies of both diagnostic tests is not rejected, and if the global test is significant then the  
 234 causes of the significance will be investigated: 1) testing the tests  $H_0 : Se_1 = Se_2$  and  
 235  $H_0 : Sp_1 = Sp_2$  individually applying the Wald test (or the likelihood ratio test) along with the

236 Holm method (or Bonferroni) to an  $\alpha$  error if the sample size is  $\leq 100$  or if the sample size is  
 237  $\geq 1000$ ; or 2) testing the tests  $H_0 : Se_1 = Se_2$  and  $H_0 : Sp_1 = Sp_2$  individually applying the  
 238 McNemar test with continuity correction ( $cc$ ) to an  $\alpha$  error if  $100 < n < 1000$ . McNemar test  
 239 statistics with  $cc$  are

$$240 \quad \chi^2_{MccSe} = \frac{(|s_{10} - s_{01}| - 1)^2}{s_{10} + s_{01}} \text{ and } \chi^2_{MccSp} = \frac{(|r_{10} - r_{01}| - 1)^2}{r_{10} + r_{01}},$$

241 respectively. In all of these test statistics we consider the frequencies of discordant pairs  $s_{ij}$   
 242 and  $r_{ij}$  with  $i \neq j$ , which are the base of the development of the McNemar test.

243 Regarding the confidence intervals for the difference between the two sensitivities  
 244 (specificities), these consist of intervals for the difference between the two paired binomial  
 245 proportions. Fagerland et al [8] compared different intervals and recommended using the  
 246 Wald interval with Bonett-Laplace adjustment. For the difference between the two  
 247 sensitivities, the Wald interval with Bonett-Laplace adjustment is

$$248 \quad Se_1 - Se_2 \in \frac{s_{10} - s_{01}}{s+2} \pm z_{1-\alpha/2} \sqrt{\frac{s_{10} + s_{01} + 2}{(s+2)^2} - \frac{(s_{10} - s_{01})^2}{(s+2)^3}},$$

249 and for the difference between the two specificities the confidence interval is

$$250 \quad Sp_1 - Sp_2 \in \frac{r_{10} - r_{01}}{r+2} \pm z_{1-\alpha/2} \sqrt{\frac{r_{10} + r_{01} + 2}{(r+2)^2} - \frac{(r_{10} - r_{01})^2}{(r+2)^3}}.$$

251 These intervals are included in the interval [-1, 1].

252 The “compbdt” program uses the method of Roldán-Nofuentes and Sidaty-Regad [6] and  
 253 the confidence interval of Wald interval with Bonett-Laplace adjustment for the difference  
 254 between the two sensitivities (specificities).

255

256

257

258    *Comparison of the likelihood ratios*

259    The comparison of the *LRs* of two diagnostic tests subject to a paired design has been the  
260    subject of several studies. Leisenring and Pepe [8] have studied the estimation of the *LRs* of a  
261    diagnostic test using a regression model, and Pepe [1] has adapted this model to compare the  
262    *LRs* individually of two binary diagnostic tests, i.e. to solve the tests  $H_0 : PLR_1 = PLR_2$  vs  
263     $H_1 : PLR_1 \neq PLR_2$  and  $H_0 : NLR_1 = NLR_2$  vs  $H_1 : NLR_1 \neq NLR_2$ . Roldán-Nofuentes and Luna  
264    [10] have compared the *LRs* individually, and also simultaneously (i.e. solving the global  
265    hypothesis                      test                       $H_0 : (PLR_1 = PLR_2 \text{ and } NLR_1 = NLR_2)$                       vs  
266     $H_1 : (NLR_1 \neq NLR_2 \text{ and/or } NLR_1 \neq NLR_2)$ , applying the maximum likelihood method. Dolgun  
267    et al [11] have extended the method of Leisenring and Pepe to compare the *LRs*  
268    simultaneously. The test statistics of the individual hypotheses tests of Pepe and the test  
269    statistics of the individual hypotheses tests of Roldán-Nofuentes and Luna have a very similar  
270    asymptotic behaviour. The test statistic of the global hypothesis test of Dolgun et al and the  
271    test statistic of the global hypothesis test of Roldán-Nofuentes and Luna have a very similar  
272    asymptotic behaviour. Therefore, the “compbdt” uses the tests proposed by Roldán-Nofuentes  
273    and Luna.

274    The method of Roldán-Nofuentes and Luna [10] compares the *LRs* considering the  
275    Napierian logarithm of the ratios of the *PLRs* and of the *NLRs*. The test statistic for the global  
276    hypothesis test of simultaneous comparison of the *LRs* is obtained applying the Wald test, i.e.

$$\chi^2_W = \hat{\boldsymbol{\omega}}^T \Sigma_{\hat{\boldsymbol{\omega}}}^{-1} \hat{\boldsymbol{\omega}},$$

278    and whose distribution is a chi-square with two degrees of freedom when the null hypothesis  
279    is true, where  $\hat{\boldsymbol{\omega}} = \left( \ln\left(\hat{PLR}_1 / \hat{PLR}_2\right), \ln\left(\hat{NLR}_1 / \hat{NLR}_2\right) \right)^T$  and  $\Sigma_{\hat{\boldsymbol{\omega}}}$  it is the estimated variance-  
280    covariance matrix obtained by applying the delta method. Roldán-Nofuentes and Amro [12]  
281    proposed the following procedure to compare the *LRs*: 1) Solve the global hypothesis test to

282 an  $\alpha$  error calculating the Wald test statistic; 2) If the global hypothesis test is not significant  
 283 to an  $\alpha$  error, then the homogeneity of the  $LRs$  of the two diagnostic tests is not rejected, but  
 284 if the global hypothesis test is significant to an  $\alpha$  error, then the study of the causes of the  
 285 significance is performed by solving the two individual hypothesis tests along with a multiple  
 286 comparison method (e.g. Holm method) to an  $\alpha$  error. In this situation, the test statistic for  
 287 comparing the two  $PLRs$  is

$$288 \quad \frac{\ln(\hat{PLR}_1/\hat{PLR}_2)}{\sqrt{\hat{Var}[\ln(\hat{PLR}_1/\hat{PLR}_2)]}}$$

289 and the test statistic for comparing the two  $NLRs$  is

$$290 \quad \frac{\ln(\hat{NLR}_1/\hat{NLR}_2)}{\sqrt{\hat{Var}[\ln(\hat{NLR}_1/\hat{NLR}_2)]}}.$$

291 These test statistics are distributed asymptotically according to a standard normal distribution.

292 Regarding the confidence intervals, Roldán-Nofuentes and Sidaty-Regad [13] studied the  
 293 comparison of the  $LRs$  through confidence intervals. For the  $PLRs$ , it is recommended to use  
 294 an interval based on the Napierian logarithm of the ratio between both, and for the  $NLRs$  it is  
 295 recommended to use a Wald type interval for the ratio between both, i.e.

$$296 \quad \frac{PLR_1}{PLR_2} \in \frac{\hat{PLR}_1}{\hat{PLR}_2} \times \exp\left[\pm z_{1-\alpha/2} \sqrt{\hat{Var}(\hat{PLR}_1/\hat{PLR}_2)}\right]$$

297 and

$$298 \quad \frac{NLR_1}{NLR_2} \in \frac{\hat{NLR}_1}{\hat{NLR}_2} \times \left[1 \pm z_{1-\alpha/2} \sqrt{\hat{Var}(\hat{NLR}_1/\hat{NLR}_2)}\right],$$

299 where the variances are calculated by applying the delta method.

300

301

302

303    *Comparison of the predictive values*

304    Comparison of the *PVs* has also been the subject of different studies. Leisenring et al [14],  
305    Wang et al [15], Kosinski [16] and Tsou [17] studied asymptotic methods to compare the  
306    *PPVs* and the *NPVs* of two diagnostic tests independently, i.e. solving the two hypothesis tests  
307     $H_0 : PPV_1 = PPV_2$  and  $H_0 : NPV_1 = NPV_2$  each one of them to an  $\alpha$  error. Takahashi and  
308    Yamamoto [18] proposed an exact test to solve this same problem. The Kosinski method has  
309    a better asymptotic performance (in terms of type I error and power) than the methods of  
310    Leisenring et al and of Wang et al. The method of Tsou leads to the same results as the  
311    Kosinski method. The method of Takahashi and Yamamoto is very conservative (as it is an  
312    exact test), even more so than the Kosinski method with small samples. The test statistics of  
313    the Kosinski method for  $H_0 : PPV_1 = PPV_2$  is

$$314 \quad T_{VPP}^{WGS} = \frac{(\hat{PPV}_1 - \hat{PPV}_2)^2}{\{\hat{PPV}_p(1 - \hat{PPV}_p) - 2C_p^{PPV}\} \left( \frac{1}{n_{10} + n_{11}} + \frac{1}{n_{01} + n_{11}} \right)}$$

315    and the test statistics for  $H_0 : NPV_1 = NPV_2$  is

$$316 \quad T_{VPN}^{WGS} = \frac{(\hat{NPV}_1 - \hat{NPV}_2)^2}{\{\hat{NPV}_p(1 - \hat{NPV}_p) - 2C_p^{NPV}\} \left( \frac{1}{n_{00} + n_{01}} + \frac{1}{n_{00} + n_{10}} \right)},$$

317    and where

$$318 \quad \hat{PPV}_p = \frac{2s_{11} + s_{10} + s_{01}}{2n_{11} + n_{10} + n_{01}}, \quad \hat{NPV}_p = \frac{2r_{00} + r_{01} + r_{10}}{2n_{00} + n_{01} + n_{10}},$$

$$319 \quad C_p^{PPV} = \frac{s_{11}(1 - \hat{PPV}_p)^2 + r_{11}\hat{PPV}_p^2}{2n_{11} + n_{10} + n_{01}} \quad \text{and} \quad C_p^{NPV} = \frac{s_{00}\hat{NPV}_p^2 + r_{00}(1 - \hat{NPV}_p)^2}{2n_{00} + n_{01} + n_{10}}.$$

320    Each statistic is distributed according to a chi-square distribution with one degree of freedom  
321    when the corresponding null hypothesis is true.

322 Roldán-Nofuentes et al [19] demonstrated that the comparison of the *PVs* of two diagnostic  
 323 tests subject to a paired design should be carried out simultaneously, i.e. solving the  
 324 hypothesis test

$$325 \quad H_0 : (PPV_1 = PPV_2 \text{ and } NPV_1 = NPV_2) \text{ vs } H_1 : (PPV_1 \neq PPV_2 \text{ and } NPV_1 \neq NPV_2).$$

326 Roldán-Nofuentes et al deduced a statistic applying the Wald test, whose distribution is a chi-  
 327 square with two degrees of freedom when the null hypothesis is true. This test statistic is

$$328 \quad \chi_w^2 = \hat{\eta}^T \Phi^T (\Phi \hat{\Sigma} \Phi^T)^{-1} \Phi \hat{\eta},$$

329 where  $\hat{\eta} = (\hat{PPV}_1, \hat{PPV}_2, \hat{NPV}_1, \hat{NPV}_2)^T$ ,  $\hat{\Sigma}$  is the estimated variance-covariance matrix of  $\hat{\eta}$   
 330 calculated by applying the delta method and  $\Phi$  is the design matrix, i.e.

$$331 \quad \Phi = \begin{pmatrix} 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & -1 \end{pmatrix}.$$

332 The test statistic  $\chi_w^2$  is distributed asymptotically according to a central chi-square  
 333 distribution with two degrees of freedom if  $H_0$  is true. Setting an  $\alpha$  error, if the global test is  
 334 not significant then we do not reject the equality of the *PVs* of both diagnostic tests; if the  
 335 global test is significant, then the investigation of the causes of the significance is carried out  
 336 applying an individual test along with a multiple comparison method (e.g. the Holm method)  
 337 to an  $\alpha$  error. The program uses the method of Roldán-Nofuentes et al [19], and as an  
 338 individual method the Kosinski method is used (calculating the weighted generalized score  
 339 statistic) since its performance is better than that of the rest of the methods.

340 Regarding the confidence intervals for the difference between the two *PPVs* and between  
 341 the two *NPVs*, these are obtained inverting the statistic of the Kosinski method, i.e.

$$342 \quad PPV_1 - PPV_2 \in \hat{PPV}_1 - \hat{PPV}_2 \pm z_{1-\alpha/2} \sqrt{\left\{ \hat{PPV}_p \left( 1 - \hat{PPV}_p \right) - 2C_p^{PPV} \right\} \left( \frac{1}{n_{10} + n_{11}} + \frac{1}{n_{01} + n_{11}} \right)}$$

343 and

344       $NPV_1 - NPV_2 \in \hat{NPV}_1 - \hat{NPV}_2 \pm z_{1-\alpha/2} \sqrt{\left\{ \hat{NPV}_p \left( 1 - \hat{NPV}_p \right) - 2C_p^{NPV} \right\} \left( \frac{1}{n_{00} + n_{01}} + \frac{1}{n_{00} + n_{10}} \right)}.$

345

346    **The “compbdt” program**

347    The “compbdt” program is a program written with R software [3] which allows us to estimate  
 348    and compare the previous parameters of two diagnostic test. The program is run with the  
 349    command

350                 $\text{compbdt}(s_{11}, s_{10}, s_{01}, s_{00}, r_{11}, r_{10}, r_{01}, r_{00})$

351    when  $\alpha = 5\%$ , and with the command

352                 $\text{compbdt}(s_{11}, s_{10}, s_{01}, s_{00}, r_{11}, r_{10}, r_{01}, r_{00}, \alpha)$

353    when  $\alpha \neq 5\%$ . Firstly, the program checks that the values introduced are viable (i.e., that  
 354    there are no negative values, values of frequencies with decimals, etc...) and that the  
 355    estimated Youden index of each diagnostic test is greater than 0 (a necessary condition for  
 356    every binary diagnostic test). The program also checks that it is possible to estimate and  
 357    compare all of the parameters. If this is not possible (for example, when there are too many  
 358    frequencies equal to 0), the program provides a message alerting to the error or the  
 359    impossibility of estimating or comparing the parameters. By default, the program shows the  
 360    numerical results with three decimal figures, a number which may be modified changing the  
 361    command “decip <- 3” at the start of the code of the program.

362    Once it is established that it is possible to carry out the study, firstly the disease prevalence  
 363    is estimated and we then estimate and compare the sensitivities and specificities, the  
 364    likelihood ratios and the predictive values, following the methods described in the previous  
 365    Section. For each type of parameter (*Se* and *Sp*, *PLR* and *NLR*, *PPV* and *NPV*), we calculate  
 366    its estimation, standard error and confidence interval to  $100(1-\alpha)\%$ . Regarding the  
 367    comparisons, if the global hypothesis test is significant, then the program solves the

368 individual hypothesis tests along with the Holm method (which is a less conservative method  
369 than the Bonferroni method) to a set  $\alpha$  error. For the hypothesis tests which are declared  
370 significant, the confidence intervals are calculated for the difference (or ratio) of the  
371 parameters. These intervals are always calculated in such a way that they are positive (for the  
372 sensitivities, specificities and predictive values), and higher than 1 for the *LRs*, indicating the  
373 diagnostic test (*Test 1* or *Test 2*) for which the parameter is estimated to be greater. If the  
374 global hypothesis test is not rejected, then the homogeneity of the parameters of both  
375 diagnostic tests is not rejected. In this situation, we do not calculate the confidence intervals  
376 for the difference or ratio of the parameters (since the homogeneity of the parameters is not  
377 rejected).

378 Furthermore, when the null hypothesis of the global hypothesis test is not rejected (and as  
379 long as the estimations are different), the program estimates the probability of making a type  
380 II error through Monte Carlo simulations. For this purpose, the program generates 10,000  
381 random samples of a multinomial distribution with the same size as the original sample and as  
382 probabilities the relative frequencies observed in the original sample. The random samples are  
383 generated in such a way that in all of them it is possible to estimate the parameters and apply  
384 the hypothesis tests. Therefore, if for one generated sample it is not possible to apply a  
385 hypothesis test, then another sample is generated instead until completing the 10,000 samples.  
386 The estimation of the probability of making a type II error is based on the data observed in the  
387 original sample i.e. the probability of making a type II error is estimated assuming that subject  
388 to the alternative hypothesis the aim is to find a difference between the parameters such as the  
389 one observed in the original sample. The estimation of this probability is of great use for  
390 researchers as the non-rejection of the null hypothesis with a probability of making a type II  
391 error greater than 20% (a value which is normally considered to be a maximum value for this  
392 probability) indicates that the null hypothesis is not reliable, and it is necessary to increase the

393 sample size. If in each global hypothesis test the alternative hypothesis test is accepted, then  
394 the program shows the estimated power of the test (one less the probability of making a type  
395 II error).

396 The results obtained comparing the sensitivities and specificities are recorded in the file  
397 “Results\_Comparison\_Accuracies.txt”, those obtained when comparing the *LRs* are recorded  
398 in the file “Results\_Comparison\_LRs.txt”, and those obtained when comparing the *PVs* are  
399 recorded in the file “Results\_Comparison\_PVs.txt”.

400

## 401 **Results**

402 The “compbdt” program has been applied to the study of Weiner et al [20] on the diagnosis of  
403 coronary artery disease, which is a classic example to illustrate statistical methods to compare  
404 parameters of two diagnostic tests. Weiner et al [19] studied the diagnosis of coronary artery  
405 disease (*CAD*) using as diagnostic tests the exercise test (*Test 1*) and the clinical history of  
406 chest pain (*Test 2*), and the coronary angiography as the gold standard. Table 2 shows the  
407 frequencies obtained by applying three medical tests to a sample of 871 individuals.

408

409 **Table 2. Study of Weiner et al.**

	Test 1 positive		Test 1 negative		Total
	Test 2 positive	Test 2 negative	Test 2 positive	Test 2 negative	
CAD	473	29	81	25	608
No CAD	22	46	44	151	263
Total	495	75	125	176	871

410

411 Running the “compbdt” program with the command

412  $\text{compbdt}(473, 29, 81, 25, 22, 46, 44, 151)$

413 the following results are obtained:

414

415

416 PREVALENCE OF THE DISEASE

417 Estimated prevalence of the disease is 69.805 % and its standard error is 0.016

418 95% confidence interval for the prevalence of the disease is (66.681% ; 72.768%)

419

420 COMPARISON OF THE ACCURACIES (SENSITIVITIES AND SPECIFICITIES)

421 Estimated sensitivity of Test 1 is 82.566% and its standard error is 0.015

422 95% confidence interval for the sensitivity of Test 1 is (79.363% ; 85.389%)

423 Estimated sensitivity of Test 2 is 91.118% and its standard error is 0.012

424 95% confidence interval for the sensitivity of Test 1 is (88.61% ; 93.148%)

425 Estimated specificity of Test 1 is 74.144% and its standard error is 0.027

426 95 % confidence interval for the specificity of Test 1 is (68.557% ; 79.087%)

427 Estimated specificity of Test 2 is 74.905% and its standard error is 0.027

428 95% confidence interval for the specificity of Test 1 is (69.358% ; 79.787%)

429 Wald test statistic for the global hypothesis test  $H_0: (Se_1 = Se_2 \text{ and } Sp_1 = Sp_2)$  is 25.662.

430 Global p-value is 0. Applying the global Wald test (to an alpha error of 5%), we reject the

431 hypothesis  $H_0: (Se_1 = Se_2 \text{ and } Sp_1 = Sp_2)$ . Estimated power (to an alpha error of 5 %) is

432 99.8%. Investigation of the causes of significance:

433 McNemar test statistic (with cc) for  $H_0: Se_1 = Se_2$  is 23.645 and the two-sided p-value is 0

434 McNemar test statistic (with cc) for  $H_0: Sp_1 = Sp_2$  is 0.011 and the two-sided p-value is

435 0.991

436 Applying the Holm method (to an alpha error of 5%), we reject the hypothesis  $H_0: Se_1 =$

437  $Se_2$  and we do not reject the hypothesis  $H_0: Sp_1 = Sp_2$

438 Sensitivity of Test 2 is significantly greater than sensitivity of Test 1. 95% confidence

439 interval for the difference  $Se_2 - Se_1$  is (5.192% ; 11.857%)

440

## COMPARISON OF THE LIKELIHOOD RATIOS

442 Estimated positive LR of Test 1 is 3.193 and its standard error is 0.339

443 95% confidence interval for the positive LR of Test 1 is (2.61 ; 3.952)

444 Estimated positive LR of Test 2 is 3.631 and its standard error is 0.39

445 95% confidence interval for the positive LR of Test 1 is (2.962 ; 4.505)

446 Estimated negative LR of Test 1 is 0.235 and its standard error is 0.022

447 95% confidence interval for the negative LR of Test 1 is (0.195 ; 0.283)

448 Estimated negative LR of Test 2 is 0.119 and its standard error is 0.016

449 95 % confidence interval for the negative LR of Test 2 is (0.09 ; 0.153)

450 Test statistic for the global hypothesis test  $H_0$ : (PLR1 = PLR2 and NLR1 = NLR2) is 23.438.

451 Global p-value is 0. Applying the global hypothesis test (to an alpha error of 5 %), we reject

452 the hypothesis  $H_0$ : (PLR1 = PLR2 and NLR1 = NLR2). Estimated power (to an alpha error of

453 5%) is 99.78 %. Investigation of the causes of significance:

454 Test statistic for H0: PLR1 = PLR2 is 0.898 and the two-sided p-value is 0.369

455 Test statistic for H0: NLR1 = NLR2 is 4.663 and the two-sided p-value is 0

456 Applying the Holm method (to an alpha error of 5 %), we do not reject the hypothesis H<sub>0</sub>:

457 PLR1 = PLR2 and we reject the hypothesis H0: NLR1 = NLR2

458 Negative likelihood ratio of Test 1 is significantly greater than negative likelihood ratio of

459 Test 2. 95% confidence interval for the ratio NLR1 / NLR2 is (1.412 ; 2.554)

460

## 461 COMPARISON OF THE PREDICTIVE VALUES

462 Estimated positive PV of Test 1 is 88.07% and its standard error is 0.014

463 95% confidence interval for the positive PV of Test 1 is (85.17% ; 90.4%

464    Estimated positive PV of Test 2 is 89.355% and its standard error is 0.012

466 Estimated negative PV of Test 1 is 64.784% and its standard error is 0.028  
467 95% confidence interval for the negative PV of Test 1 is (59.246% ; 69.976%)  
468 Estimated negative PV of Test 2 is 78.486% and its standard error is 0.026  
469 95% confidence interval for the negative PV of Test 2 is (73.024% ; 83.151%)  
470 Wald test statistic for the global hypothesis test  $H_0$ : ( $PPV_1 = PPV_2$  and  $NPV_1 = NPV_2$ ) is  
471 25.944. Global p-value is 0. Applying the global hypothesis test (to an alpha error of 5%), we  
472 reject the hypothesis  $H_0$ : ( $PPV_1 = PPV_2$  and  $NPV_1 = NPV_2$ ). Estimated power (to an alpha  
473 error of 5%) is 99.26%. Investigation of the causes of significance:  
474 Weighted generalized score statistic for  $H_0$ :  $PPV_1 = PPV_2$  is 0.807 and the two-sided p-  
475 value is 0.369  
476 Weighted generalized score statistic for  $H_0$ :  $NPV_1 = NPV_2$  is 22.502 and the two-sided p-  
477 value is 0  
478 Applying the Holm method (to an alpha error of 5%), we do not reject the hypothesis  $H_0$ :  
479  $PPV_1 = PPV_2$  and we reject the hypothesis  $H_0$ :  $NPV_1 = NPV_2$   
480 Negative PV of Test 2 is significantly greater than negative PV of Test 1. 95% confidence  
481 interval for the difference  $NPV_2 - NPV_1$  is (8.041% ; 19.363%)  
482  
483 These outputs obtained when running the program allow researchers to interpret the results  
484 easily. First, for each type of parameters, all parameters are estimated and the corresponding  
485 global test is solved. In summary, the three global hypothesis tests are rejected and then the  
486 causes of the significance of each global test are investigated. For individual hypothesis tests  
487 that are declared significant, it is indicated which is the diagnostic test for which the  
488 parameter is greater, calculating the corresponding confidence interval. Due to the high  
489 sample size, the estimated power for each of the global tests is very high (close to 100%).

490 In R, an alternative program to “compbdt” is the DTComPair package [21]. The  
 491 DTComPair package estimates the same parameters as the “compbdt” and compares the  
 492 parameters individually, i.e. solving each hypothesis test to an  $\alpha$  error. Table 3 shows the  
 493 results obtained when applying the DTComPair package with  $\alpha = 5\%$  (the estimations of the  
 494 parameters and their standard errors are not shown as they are the same as those obtained with  
 495 the “compbdt” program). The conclusions obtained are similar to those obtained with the  
 496 “compbdt” program, although this program uses methods with better asymptotic behaviour.  
 497

498 **Table 3. Results obtained with the DTComPair package.**

Confidence intervals for the parameters of each diagnostic test (95% confidence)		
	Test 1	Test 2
Sensitivity	79.550% ; 85.582%	88.857% ; 93.380%
Specificity	68.853% ; 79.436%	69.665% ; 80.145%
Positive LR	2.594 ; 3.931	2.942 ; 4.481
Negative LR	0.195 ; 0.284	0.091 ; 0.154
Positive PV	85.409% ; 90.731%	86.927% ; 91.783%
Negative PV	59.388% ; 70.180%	73.403% ; 83.570%
Comparison of the parameters of the two diagnostic tests ( $\alpha = 5\%$ )		
Sensitivities		
McNemar test statistics: test statistic = 24.582, p-value = 0		
Exact test: p-value = 0		
95% Tango confidence interval for $Se_2 - Se_1$ : 5.278% ; 11.966		
Specificities		
McNemar test: test statistic = 0.044, p-value = 0.833		
Exact test: two-sided p-value = 0.916		
Likelihood ratios (Method of Leisenring et al [9] and Pepe [1])		
Positive LRs: test statistic = -0.898, p-value = 0.369		
Negative LRs: test statistic = 4.663, p-value = 0		
95% confidence interval for $NLR_1/NLR_2$ : 1.487 ; 2.644		
Predictive values (Method of Leisenring et al [14])		
Positive PVs: test statistic = 0.802, p-value = 0.371		
Negative PVs: test statistic = 23.579, p-value = 0		
Predictive values (Method of Kosinski [15])		
Positive PVs: test statistic = 0.807, p-value = 0.369		
Negative PVs: test statistic = 22.502, p-value = 0		
Relative predictive values (Method of Moskowitz and Pepe [26])		
Positive PVs: test statistic = -0.895, p-value = 0.371		
Negative PVs: test statistic = -4.737, p-value = 0		
95% confidence interval for $NPV_1/NPV_2$ : 0.762 ; 0.894		

499

500 **Conclusions**

501 The comparison of the performance of two diagnostic tests subject to a paired design is an  
502 important topic in Medicine. Many studies have been carried out on statistical methods to  
503 estimate and compare parameters of two binary diagnostic tests subject to this type of design.  
504 In the “compbdt” program the most efficient methods have been implemented, in terms of  
505 coverage and width for the confidence intervals and in terms of type I error and power for the  
506 hypothesis tests, developed up to the present day. The comparisons of the three types of  
507 parameters (sensitivities and specificities, likelihood ratios and predictive values) are based on  
508 solving the global hypothesis tests. For each type of parameter, the program solves the global  
509 test and if this is not significant to an  $\alpha$  error then we do not reject the homogeneity of the  
510 parameters of both diagnostic tests; if the global test is significant to an  $\alpha$  error then the  
511 causes of the significance are investigated solving the individual hypothesis tests along with  
512 Holm’s method of multiple comparison to an  $\alpha$  error. This procedure is very similar to  
513 analysis of variance. If for each type of parameter we directly solve each one of the individual  
514 hypothesis tests to an  $\alpha$  error, it is possible to obtain mistaken results. Two examples of this  
515 are explained in the articles by Roldán-Nofuentes and Sidaty-Regad [6] and Roldán-  
516 Nofuentes et al [19].

517 The program requires installing the *R* software, which is freely available at the URL  
518 “<https://www.r-project.org>”, and it is necessary for the data observed to have the structure  
519 given in Table 1. The program provides all of the results necessary so that the researcher can  
520 make interpretations in a simple way. Another contribution made by this program is the  
521 estimation of the probability of making a type II error based on the data observed in the  
522 sample through Monte Carlo simulations, data which provides information about the  
523 reliability of the null hypothesis when the hypothesis test is not significant. The program has

524 been applied to a classic example of this topic. On an Intel Core i7 3.40 GHz computer the  
525 program has been run in around 7 seconds.

526 With respect to the DTComPair package [21], the “compbdt” program uses methods with  
527 better asymptotic behaviour and has the following advantages:

528 a) For a binomial proportion (such as the sensitivity, specificity and predictive values of  
529 each diagnostic test), the DTComPair package uses the Agresti and Coull interval [22]. The  
530 “compbdt” uses the interval of Yu et al [4], which has a better coverage than that of Agresti  
531 and Coull.

532 b) The DTComPair uses the interval of Simel et al [23] for the positive (negative)  
533 likelihood ratio of each diagnostic test, an interval which, as is well known, does not have a  
534 good coverage when the samples are not very large. The “compbdt” program uses the interval  
535 of Martín-Andrés and Álvarez-Hernández, which is the interval with the best coverage for the  
536 ratio of two independent binomial proportions (such as the positive and negative likelihood  
537 ratios).

538 c) The DTComPair package compares the parameters individually, which can lead to  
539 mistakes [6, 19]. The “compbdt” program is based on the simultaneous comparisons of the  
540 parameters and on research into the causes of the significance when the global tests are  
541 significant.

542 d) The DTComPair package calculates three confidence intervals for the difference of the  
543 two sensitivities (specificities): Wald (with or without  $cc$ ), Agresti and Min [24], and Tango  
544 [25]. Fagerland et al [8] have shown that the Wald interval with Bonett-Laplace adjustment  
545 (interval implemented in the “compbdt” program) has an asymptotic behaviour very similar to  
546 that of Tango, and that both intervals have a better behaviour than that of Agresti and Min.  
547 The advantage of the Wald interval with Bonett-Laplace adjustment is that this interval has  
548 closed-form expression.

549 e) The DTComPair package calculates confidence intervals for the ratio of *LRs* based on  
550 regression models [1, 9]. The “compbdt” program uses confidence intervals with better  
551 asymptotic behaviour [13].

552 f). The “compbdt” program estimates the power or probability of making a type II error,  
553 depending on whether or not the alternative hypothesis is accepted or not the null hypothesis  
554 is rejected, based on the data observed in the sample through Monte Carlo simulations.

555 g) The DTComPair package only provides numerical results, whereas the “compbdt”  
556 program also interprets them, which is of great use for the clinician.

557 The application of the “compbdt” program requires the results of both diagnostic tests and  
558 the gold standard to be known for all of the individuals in the sample. If the result of a  
559 diagnostic test is unknown for any individual, and this missing data is random due to chance  
560 (the missing data mechanism is missing at random), this data can always be imputed applying  
561 some method of imputation and then it is possible to use the program to solve the problem of  
562 comparison of the parameters. The program also requires knowledge of the discordant  
563 frequencies ( $s_{ij}$  and  $r_{ij}$  with  $i \neq j$ ), since these are necessary to be able to solve the  
564 hypothesis tests. If the researcher wants to use the “compbdt” program to repeat the results of  
565 a study and we do not know the discordant frequencies but we do know an estimation of the  
566 Cohen kappa coefficient (or another measure of association) between the diagnostic tests in  
567 diseased individuals and in non-diseased individuals, then it is possible to use both  
568 estimations to obtain the values of the discordant frequencies. The “compbdt” program is  
569 available as supplementary material of this manuscript.

570 Finally, the “compbdt” program can also be applied when the sampling is case-control, i.e.  
571 the two diagnostic tests are applied to two samples, one of  $n_1$  diseased individuals and  
572 another one of  $n_2$  non-diseased individuals. In this situation, the frequencies  $s_{ij}$  correspond to

573 the case sample (with  $n_1 = \sum_{i,j=0}^1 s_{ij}$ ) and the frequencies  $r_{ij}$  correspond to the control sample  
574 (with  $n_2 = \sum_{i,j=0}^1 r_{ij}$ ). Subject to this sampling, it is necessary to take into account the fact that  
575 the results obtained for the prevalence and all of the results obtained for the predictive values  
576 are not valid, since from a case-control sample it is not possible to obtain an estimation of the  
577 disease prevalence (the value  $n_1/(n_1+n_2)$  is not an estimation of the prevalence since the  
578 sample sizes  $n_1$  and  $n_2$  are set by the researcher).

579

## 580 **Abbreviations**

581 CAD: Coronary Artery Disease; LR: Likelihood ratio; p: Prevalence of the disease; PLR:  
582 Positive Likelihood Ratio; PV: Predictive Value; PPV: Positive Predictive Value; NLR:  
583 Negative Likelihood Ratio NPV: Negative Predictive Value; Se: Sensitivity; Sp: Specificity.

584

## 585 **Acknowledgements**

586 I thank the Editor and the two referees for their helpful comments that improved the quality of  
587 the manuscript.

588

## 589 **Availability and requirements**

590 Project name: Comparison of binary diagnostic tests  
591 Project home page: <https://www.ugr.es/~bioest/>  
592 Operating system(s): Platform independent  
593 Programming language: R  
594 Other requirements: R 3.6.1 or above  
595 License: GPL-2

596 Any restrictions to use by non-academics: none

597

598 **Ethics approval and consent to participate**

599 Not applicable.

600

601 **Consent for publication**

602 Not applicable.

603

604 **Availability of data and material**

605 All data generated or analyzed during the current study are included in this published article.

606 The program “compbdt” is available as supplementary material of this manuscript.

607

608 **Competing interests**

609 The author declares that they have no conflict of interest.

610

611 **Funding**

612 This research was supported by the Spanish Ministry of Economy, Grant Number MTM2016-

613 76938-P.

614

615 **Author's contributions**

616 RN has reviewed the statistical methods and has written the program. The author reads and

617 approves the final manuscript.

618

619

620

621 **References**

- 622 1 Pepe, MS. The statistical evaluation of medical tests for classification and prediction. New  
623 York: Oxford University Press; 2003.
- 624 2 Zhou, XH, Obuchowski, NA. McClish, DK. Statistical Methods in Diagnostic  
625 Medicine, second edition. New York: Wiley; 2011.
- 626 3 R.C. Team R. A Language and Environment for Statistical Computing. Vienna, Austria,  
627 2016, URL <https://www.R-project.org/>
- 628 4 Yu, W, Guo, X, Xu, W. An improved score interval with a modified midpoint for a  
629 binomial proportion, Journal of Statistical Computation and Simulation, 2014;84:1022-  
630 1038.
- 631 5 Martín-Andrés, A, Álvarez-Hernández, M. Two-tailed approximate confidence intervals  
632 for the ratio of proportions. Statistics and Computing, 2014;24:65-75.
- 633 6 Roldán-Nofuentes, JA, Sidaty-Regad, SB. Recommended methods to compare the  
634 accuracy of two binary diagnostic tests subject to a paired design. Journal of Statistical  
635 Computation and Simulation, 2019;89:2621-2644.
- 636 7 Holm, S. A simple sequential rejective multiple testing procedure. Scandinavian Journal  
637 of Statistics. 1979;6:65-70.
- 638 8 Fagerland, MW, Lydersen, S, Laake, P. Recommended tests and confidence intervals  
639 for paired binomial proportions. Statistics in Medicine, 2014;33:2850-2875.
- 640 9 Leisenring, W, Pepe, MS. Regression modelling of diagnostic likelihood ratios for the  
641 evaluation of medical diagnostic tests. Biometrics, 1998;54:444-442.
- 642 10 Roldán-Nofuentes, JA, Luna del Castillo, JD. Comparison of the likelihood ratios of two  
643 binary diagnostic tests in paired designs. Statistics in Medicine, 2007;26:4179-4201.

- 644 11 Dolgun, NA, Gozukara, H, Karaagaoglu, E. Comparing diagnostic tests: test of  
645 hypothesis for likelihood ratios. *Journal of Statistical Computation and Simulation*,  
646 2012;82:369-381.
- 647 12 Roldán-Nofuentes, JA, Amro, R. Estimation and comparison of the likelihood ratios of  
648 binary diagnostic tests. In: M. Negreiros M, Bouza C, Mello F, editors. *Models and*  
649 *methods for supporting decision making in human health and environment protection*.  
650 New York: Nova Science Publishers; 2016. p. 57-70.
- 651 13 Roldán-Nofuentes, JA, Sidaty-Regad, SB. Comparison of the likelihood ratios of two  
652 diagnostic tests subject to a paired design: confidence intervals and sample size. *Revstat*  
653 *Statistical Journal*, 2020;in press.
- 654 14 Leisenring, W, Alonzo, T, Pepe, MS. Comparisons of predictive values of binary  
655 medical diagnostic tests for paired designs. *Biometrics*, 2000;56:345-351.
- 656 15 Wang, W, Davis, CS, Soong, SJ. Comparison of predictive values of two diagnostic  
657 tests from the same sample of subjects using weighted least squares. *Statistics in*  
658 *Medicine*, 2006;25:2215-2229.
- 659 16 Kosinski, AK. A weighted generalized score statistic for comparison of predictive  
660 values of diagnostic tests. *Statistics in Medicine*, 2013;32:964-77.
- 661 17 Tsou, TS. A new likelihood approach to inference about predictive values of diagnostic  
662 tests in paired designs. *Statistical Methods in Medical Research*, 2018;27:541-548.
- 663 18 Takahashi, K, Yamamoto, K. An exact test for comparing two predictive values in  
664 small-size clinical trials. *Pharmaceutical Statistics*, 2019;in press.
- 665 19 Roldán-Nofuentes, JA, Luna del Castillo, JD, Montero-Alonso, MA. Global hypothesis  
666 test to simultaneously compare the predictive values of two binary diagnostic tests.  
667 *Computational Statistics and Data Analysis*, 52012;6:1161-1173.

- 668 20 Weiner, DA, Ryan, TJ, McCabe, CH, Kennedy, JW, Schloss, M, Tristani, F, Chaitman,  
669 BR, Fisher, LD. Correlations among history of angina, ST-segment and prevalence of  
670 coronary artery disease in the coronary artery surgery study (CASS). New England  
671 Journal of Medicine, 1979;301:230-235.
- 672 21. Stock, C, Hielscher, T. DTComPair: comparison of binary diagnostic tests in a paired  
673 study design. R package version 1.0.3. URL: [http://CRAN.R-  
674 project.org/package=DTComPair](http://CRAN.R-project.org/package=DTComPair), 2014.
- 675 22. Agresti, A, Coull, .
- 676 23. Simel, DL, Samsa, GP, Matchar, DB. Likelihood ratios with confidence: sample size  
677 estimation for diagnostic test studies. Journal of Clinical Epidemiology, 1991;44:763-  
678 770.
- 679 24. Agresti A, Min Y. Effects and non-effects of paired identical observations in comparing  
680 proportions with binary matched-pairs data. Statistics in Medicine, 2004;23:65-75.
- 681 25. Tango T. Equivalence test and confidence interval for the difference in proportions for  
682 the paired-sample design. Statistics in Medicine, 1998;17:891-908.
- 683 26. Moskowitz, CS, Pepe, MS. Comparing the predictive values of diagnostic tests: sample  
684 size and analysis for paired study designs. Clinical Trials, 2006;3:272-279.