

Two Dependent Diagnostic Tests: Use of Copula Functions in the Estimation of the Prevalence and Performance Test Parameters

Dos pruebas para diagnóstico clínico: uso de funciones copula en la estimación de la prevalencia y los parámetros de desempeño de las pruebas

JOSÉ RAFAEL TOVAR^{1,a}, JORGE ALBERTO ACHCAR^{2,b}

¹CENTRO DE INVESTIGACIONES EN CIENCIAS DE LA SALUD (CISC), ESCUELA DE MEDICINA Y CIENCIAS DE LA SALUD, UNIVERSIDAD DEL ROSARIO, BOGOTÁ, COLOMBIA

²DEPARTAMENTO DE MEDICINA SOCIAL FMRP, FACULDADE DE SAÚDE, UNIVERSIDADE DE SÃO PAULO, RIBERÃO PRETO, BRASIL

Abstract

In this paper, we introduce a Bayesian analysis to estimate the prevalence and performance test parameters of two diagnostic tests. We concentrated our interest in studies where the individuals with negative outcomes in both tests are not verified by a gold standard. Given that the screening tests are applied in the same individual we assume dependence between test results. Generally, to capture the possible existing dependence between test outcomes, it is assumed a binary covariance structure, but in this paper, as an alternative for this modeling, we consider the use of copula function structures. The posterior summaries of interest are obtained using standard MCMC (Markov Chain Monte Carlo) methods. We compare the results obtained with our approach with those obtained using binary covariance and assuming independence. We considerate two published medical data sets to illustrate the approach.

Key words: Bayes analysis, Copula, Dependence, Monte Carlo Simulation, Public health.

Resumen

En este artículo introducimos un análisis Bayesiano para estimar la prevalencia y los parámetros de desempeño de pruebas para diagnóstico clínico, con datos obtenidos bajo estudios de tamizaje que incluyen el uso de dos pruebas diagnósticas en los cuales, los individuos con resultado negativo en

^aLecturer. E-mail: rtovar34@hotmail.com

^bAssociate Professor. E-mail: achcar@fmrp.usp.br

las dos pruebas no son confirmados con una prueba patrón de oro. Dado que las pruebas de tamizaje son aplicadas al mismo individuo, nosotros asumimos dependencia entre los resultados de las pruebas. Generalmente, para capturar la posible dependencia existente entre los resultados de las pruebas diagnósticas, se asume una estructura de covarianza binaria, pero en este artículo, nosotros consideramos el uso de estructuras que pueden ser modaladas usando funciones cópula, como una alternativa al modelamiento de la dependencia. Las estadísticas a posteriori de interés son obtenidas usando métodos MCMC. Los resultados obtenidos usando nuestra aproximación son comparados con los obtenidos usando modelos que asumen estructura binaria y con los obtenidos usando modelos bajo el supuesto de independencia entre resultados de las pruebas para diagnóstico clínico. Para ilustrar la aplicación del método y para hacer las comparaciones se usaron los datos de dos estudios publicados en la literatura.

Palabras clave: análisis bayesiano, copula, dependencia, simulación Monte Carlo, salud pública.

1. Introduction

In literature, there are designs to evaluate new screening tests in which more than one diagnostic test is applied to the same individual and where in some cases all patients cannot be verified by a test free of error to classify individuals or Gold Standard. This situation implies in the presence of verification bias. When the design considers the use of two continuous scale diagnostic tests transformed to a binary scale using a cut-off point to classify an individual as positive or negative to a given disease, these tests could have dependent outcomes within a continuous dependence structure but as we have the final binary results to do the data analysis, we could model the dependence considering a bivariate Bernoulli distribution with the covariance as a dependence parameter. This approach has been studied by different authors such as Thibodeau (1981), Vacek (1985) and Walter & Irwig (1988), amongst others. Assuming binary structure, Bohning & Patilea (2008), developed two indexes to study the dependence between two diagnostic tests: a first is derived using the λ reparametrization introduced by Georgiadis, Johnson & Gardner (2003) and a second index derived by applying the OR (odds ratio) concept on 2×2 probability tables associated with the two diagnostic test results. Some approaches such as those of Brenner (1996), Qu & Hadgu (1998) and Torrance-Rynard & Walter (1997), have considered the continuous structure in the data to study the dependence between test outcomes using models of latent variable.

In this paper, we introduce a Bayesian model to estimate the prevalence, performance test parameters and the dependence between them, using two copula functions, the FGM (Farlie-Gumbel-Morgenstern) copula and the Gumbel copula. The FGM is a copula function that allows modeling very weak linear dependencies usually not easily observed using traditional bivariate plots.

If the continuous traits that make up the diagnostic tests have a dependence like FGM structure, usually the data analyst assumes independence in the statisti-

cal model used to obtain the parameter estimates. The form of the Gumbel copula used in this work, models relatively weak negative linear dependencies but the copula parameter of dependence belongs in the space $(0,1)$. In agreement with some simulation results not showed in this paper, the bivariate plots obtained under different levels of Gumbel copula dependence show a dispersion similar with that observed when the data are obtained under independence assumption, then, it is not easy to observe the presence of a negative correlation between test outcomes. The use of this copula, also allows us to study dependencies with not necessarily linear structures which is possible in diagnostic situations whose results are obtained after dichotomization.

We compare the estimates obtained using copula models with those obtained assuming binary covariance structure and independence assumption. In our approach, we assume that the diagnostic procedure includes two (observable or not) variables measured on a continuous scale with some type of positive dependence between them that can be modeled using copula functions. Copula functions have been widely used for modeling the dependence between continuous scale variables regardless the type of distribution underlying in the margins, in many other subject or topic areas as hydrology and finance.

To illustrate our proposed models, we use two data sets introduced in the literature. The first one, was obtained from Smith, Bullock & Catalona (1997), who screened 19,476 men for prostate cancer using the Digital Rectal Exam (DRE) and the Prostate Specific Antigen (PSA) in serum. With that same data set, Bohning & Patilea (2008) and Martinez, Achcar & Louzada (2005) studied the association between diagnostic test results. The second data set was introduced by Ali, Moodambail, Hamrah, Bin-Nakhi & Sadeq (2007), where they evaluated a fast method to detect urinary tract infection in 132 children of both genders with ages ranging from three days to 11 years.

This paper is organized as follows: In Section 2 we introduce the model formulation for two associated diagnostic tests; in Section 3, we present our Bayesian estimation procedure; in Section 4, we introduce two examples; finally in section 5, we present some discussion on the obtained results.

2. Model Formulation for Two Dependent Diagnostic Tests

We consider four different models that can be used, the first model assumes conditionally independent tests results and the other three models assume that the tests are dependent conditionally on the disease status.

2.1. Model Under Independence Assumption

Two diagnostic tests are respectively denoted by T_1 and T_2 where $T_\nu = 1$ is related to a positive result for the test ν , $\nu = 1, 2$ and $T_\nu = 0$ is related to a negative result. In Table 1 we have a generic representation of the tests compared

with an ideal reference test. If the study design implies that individuals with negative outcome in both tests are not verified by a test free of error to classify the individuals (“Gold Standard”), the values d, h, n_+ and n_- (showed in brackets), are unknown although the sum $u = n_+ + n_-$ is known.

TABLE 1: Tests results. Values in brackets are unknown under verification bias.

	Diseased subjects			Non-diseased subjects		
	$T_2 = 1$	$T_2 = 0$	Total	$T_2 = 1$	$T_2 = 0$	Total
$T_1 = 1$	a	b	$a + b$	e	f	$e + f$
$T_1 = 0$	c	[d]	$c + [d]$	g	[h]	$g + [h]$
Total	$a + c$	$b + [d]$	[n_+]	$e + g$	$f + [h]$	[n_-]

Let us denote by p the prevalence of a disease and by D the true status, when $D = 1$ denotes a diseased individual and $D = 0$ denotes a non-diseased individual. That is, $p = P(D = 1)$. The sensitivities are given by $S_\nu = P(T_\nu = 1 | D = 1)$ and the specificities are given by $E_\nu = P(T_\nu = 0 | D = 0)$.

For the independence assumption model, we use the Bayesian estimation procedure developed by Martinez et al. (2005) to obtain the likelihood contributions of the eight possible combinations of results among tests and true disease state as appear in the left column in Table 2.

2.2. Model Under Binary Dependence Structure

For a binary structure model, we assume as dependence parameter, a positive covariance between tests based on the joint Bernoulli distribution. We assumed that the dependence between tests is similar in diseased and non-diseased populations in the same way as considered by Dendukuri & Joseph (2001) to obtain the contributions to likelihood function of the eight combinations of results among the two diagnostic tests and the Gold Standard. The results are showed in Table 2.

TABLE 2: Likelihood contributions of all possible combinations of outcomes of T_1, T_2 and D . (f_i = number of individuals in the cell i ; $i = 1, 2, \dots, 8$. Values in brackets are unknown under verification bias).

i	D	T_1	T_2	f_i	Contribution to likelihood	
					Independence assumption	Binary dependence
1	1	1	1	a	pS_1S_2	$p[S_1S_2 + \psi_D]$
2	1	1	0	b	$pS_1(1 - S_2)$	$p[S_1(1 - S_2) - \psi_D]$
3	1	0	1	c	$p(1 - S_1)S_2$	$p[(1 - S_1)S_2 - \psi_D]$
4	1	0	0	[d]	$p(1 - S_1)(1 - S_2)$	$p[(1 - S_1)(1 - S_2) + \psi_D]$
5	0	1	1	e	$(1 - p)(1 - E_1)(1 - E_2)$	$(1 - p)[(1 - E_1)(1 - E_2) + \psi_{ND}]$
6	0	1	0	f	$(1 - p)(1 - E_1)E_2$	$(1 - p)[(1 - E_1)E_2 - \psi_{ND}]$
7	0	0	1	g	$(1 - p)E_1(1 - E_2)$	$(1 - p)[E_1(1 - E_2) - \psi_{ND}]$
8	0	0	0	[h]	$(1 - p)E_1E_2$	$(1 - p)[E_1E_2 + \psi_{ND}]$

2.3. Model Assuming a Dependence Copula Structure

Let us assume that the test outcomes are realizations of the random variables V_1 and V_2 measured on a positive continuous scale ($V_1 > 0$ and $V_2 > 0$) which represent the expression of two biological traits whose behavior is altered by the presence of disease or infection process. Also, let us assume that two cut-off values ξ_1 and ξ_2 are chosen for each test in order to determine when an individual is classified as positive or negative. In this way we assume that an individual is classified as positive for test ν if $V_\nu > \xi_\nu$ that is, $T_\nu = 1$ if and only if $V_\nu > \xi_\nu$ for $\nu = 1, 2$. To model the dependence structure between the random variables V_1 and V_2 , let us consider the use of copula functions, which has been studied by many authors ((Nelsen 1999) is a classical book on this topic). Multivariate distribution functions F can be written in the form of a copula function, that is, if $F(v_1, \dots, v_m)$ is a joint multivariate distribution function with univariate marginal distribution functions $F_1(v_1), \dots, F_m(v_m)$, thus there exists a copula function $C(u_1, \dots, u_m)$ such that,

$$F(v_1, \dots, v_m) = C(F_1(v_1), \dots, F_m(v_m)) \quad (1)$$

When the marginal distributions are continuous, a copula function always exists and can be found from the relation

$$C(u_1, \dots, u_m) = F(F_1^{-1}(u_1), \dots, F_m^{-1}(u_m)) \quad (2)$$

For the special case of bivariate distributions, we have $m = 2$. The approach to formulate a multivariate distribution using a copula is based on the idea that a simple transformation ($U = F_1(V_1)$ and $W = F_2(V_2)$) can be made of each marginal variable in such a way that each transformed marginal variable has a uniform distribution. Specifying dependence between V_1 and V_2 is the same as specifying dependence between U and W , thus the problem reduces to specifying a bivariate distribution between two uniform variables, that is a copula.

2.3.1. Model Considering Dependence Type FGM Copula

The third model considered for the study of the dependence structure for two tests, is based in the Farlie Gumbel Morgenstern (FGM) copula widely studied by authors as Nelsen (1999), Amblard & Girard (2002, 2005, 2008). The FGM copula is defined by,

$$C_I(u, w) = uw[1 + \varphi(1 - u)(1 - w)] \quad (3)$$

where $u = F_1(v_1)$, $w = F_2(v_2)$ and φ is a copula parameter such that $-1 \leq \varphi \leq 1$. If $\varphi = 0$, we have two independent marginal random variables. We assume different parameters φ_D and φ_{ND} for diseased and non-diseased individuals, respectively.

From (3) the cumulative joint distribution and the join survival function for the random variables V_1 and V_2 is given by,

$$\begin{aligned} F_I(v_1, v_2) &= C_I(F_1(v_1), F_2(v_2)) \\ &= F_1(v_1)F_2(v_2)[1 + \varphi(1 - F_1(v_1))(1 - F_2(v_2))] \end{aligned} \quad (4)$$

$$S(v_1, v_2) = P(V_1 > v_1, V_2 > v_2) = 1 - F_1(v_1) - F_2(v_2) + F(v_1, v_2) \quad (5)$$

Within the diseased individuals group, we have,

$$\begin{aligned} F_1^D(\xi_1) &= P(V_1 \leq \xi_1 | D = 1) = 1 - S_1 \\ F_2^D(\xi_2) &= P(V_2 \leq \xi_2 | D = 1) = 1 - S_2 \end{aligned}$$

From (4), we have

$$\begin{aligned} F_D(\xi_1, \xi_2) &= F_1^D(\xi_1)F_2^D(\xi_2)[1 + \varphi(1 - F_1^D(\xi_1))(1 - F_2^D(\xi_2))] \\ &= (1 - S_1)(1 - S_2)(1 + \varphi_D S_1 S_2) \end{aligned}$$

and from (5) we have,

$$\begin{aligned} P(T_1 = 1, T_2 = 1 | D = 1) &= S_D(\xi_1, \xi_2) \\ &= 1 - (1 - S_1) - (1 - S_2) + (1 - S_1)(1 - S_2)(1 + \varphi_D S_1 S_2) \end{aligned}$$

That is,

$$P(T_1 = 1, T_2 = 1 | D = 1) = S_1 S_2 (1 + \varphi_D (1 - S_1)(1 - S_2))$$

and

$$P(T_1 = 1, T_2 = 1, D = 1) = p S_1 S_2 (1 + \varphi_D (1 - S_1)(1 - S_2))$$

Observe that, if $\varphi_D = 0$ (independent test outcomes), we have

$$P(T_1 = 1, T_2 = 1, D = 1) = p S_1 S_2$$

as given in Table 2.

Also,

$$\begin{aligned} P(T_1 = 1, T_2 = 0, D = 1) &= P(D = 1)P(T_1 = 1, T_2 = 0 | D = 1) \\ &= pP(V_1 > \xi_1, V_2 \leq \xi_2 | D = 1) \end{aligned}$$

On the other hand,

$$\begin{aligned} P(V_1 > \xi_1, V_2 \leq \xi_2 | D = 1) &= P(V_2 \leq \xi_2 | D = 1) - P(V_1 \leq \xi_1, V_2 \leq \xi_2 | D = 1) \\ &= F_2^D(\xi_2) - F_D(\xi_1, \xi_2) \end{aligned}$$

Thus,

$$P(T_1 = 1, T_2 = 1, D = 1) = p(1 - S_2)S_1[1 - \varphi_D S_2(1 - S_1)]$$

If $\varphi_D = 0$, we have

$$P(T_1 = 1, T_2 = 0, D = 1) = p S_1 (1 - S_2)$$

as in the independent case (see Table 2).

Similarly,

$$\begin{aligned} P(T_1 = 0, T_2 = 1, D = 1) &= P(D = 1)P(T_1 = 0, T_2 = 1 | D = 1) \\ &= pP(V_1 \leq \xi_1, V_2 > \xi_2 | D = 1) \end{aligned}$$

Since,

$$\begin{aligned} P(V_1 \leq \xi_1, V_2 > \xi_2 | D = 1) &= P(V_1 \leq \xi_1 | D = 1) - P(V_1 \leq \xi_1, V_2 \leq \xi_2 | D = 1) \\ &= F_1^D(\xi_1) - F_D(\xi_1, \xi_2) \end{aligned}$$

then,

$$P(T_1 = 0, T_2 = 1, D = 1) = p(1 - S_1)S_2[1 - \varphi_D S_1(1 - S_2)]$$

When $\varphi_D = 0$ we have $P(T_1 = 0, T_2 = 1, D = 1) = pS_2(1 - S_1)$ as in the independent case (see Table 2).

We also have,

$$\begin{aligned} P(T_1 = 0, T_2 = 0, D = 1) &= P(D = 1)P(T_1 = 0, T_2 = 0 | D = 1) \\ &= pP(V_1 \leq \xi_1, V_2 \leq \xi_2 | D = 1) \\ &= pF_D(\xi_1, \xi_2), \end{aligned}$$

that is,

$$P(T_1 = 0, T_2 = 0, D = 1) = p(1 - S_1)(1 - S_2)[1 + \varphi_D S_1 S_2]$$

Within the non-diseased individuals group, we have:

$$\begin{aligned} P(T_1 = 1, T_2 = 1, D = 0) &= P(D = 0)P(T_1 = 1, T_2 = 1 | D = 0) \\ &= (1 - p)P(V_1 > \xi_1, V_2 > \xi_2 | D = 0) \\ &= (1 - p)S_{ND}(\xi_1, \xi_2) \\ &= (1 - p)(1 - F_1^{ND}(\xi_1) - F_2^{ND}(\xi_2) + F_{ND}(\xi_1, \xi_2)) \end{aligned}$$

Observe that,

$$\begin{aligned} P(T_1 = 0 | D = 0) &= P(V_1 \leq \xi_1 | D = 0) = F_1^{ND}(\xi_1) = E_1 \quad \text{and} \\ P(T_2 = 0 | D = 0) &= P(V_2 \leq \xi_2 | D = 0) = F_2^{ND}(\xi_2) = E_2 \end{aligned}$$

Using (4), we have

$$F_{ND}(\xi_1, \xi_2) = E_1 E_2 [1 + \varphi_{ND}(1 - E_1)(1 - E_2)]$$

That is,

$$P(T_1 = 1, T_2 = 1, D = 0) = (1 - p)(1 - E_1)(1 - E_2)[1 + \varphi_{ND} E_1 E_2]$$

The contributions to the likelihood for all situations with diseased and non-diseased individuals are given in Table 3.

2.3.2. Model Considering Dependence Type Gumbel Copula

The last considered model, is derived from Gumbel copula function defined as;

$$C_{II}(u, w) = u + w - 1 + (1 - u)(1 - w) \exp\{-\phi \log(1 - u) \log(1 - w)\} \quad (6)$$

In this model, the joint cumulative distribution function for the random variables V_1 and V_2 is given by,

$$F_{II}(v_1 v_2) = F_1(v_1) + F_2(v_2) - 1 + (1 - F_1(v_1))(1 - F_2(v_2)) \exp\{-\phi \log(1 - F_1(v_1)) \log(1 - F_2(v_2))\} \quad (7)$$

As pointed out by (Gumbel 1960) for this copula model, when $\phi = 1$ the Pearson correlation linear coefficient (ρ) takes the value -0.40365 . In this case, the parameter of the Gumbel copula, does not models positive linear correlations. Also, when the two variables are independent, ϕ takes the zero value.

Employing the same arguments considered with the FGM copula and using (7) we obtain all the contributions for the likelihood function when it is assumed a Gumbel copula dependence structure (Table 3).

TABLE 3: Likelihood contributions of all possible combinations of outcomes of T_1, T_2 and D when the dependence has the “FGM copula” or “Gumbel copula” structure. (f_i = number of individuals in the cell i ; $i = 1, 2, \dots, 8$. Values in brackets are unknown under verification bias).

					Contribution to likelihood	
i	D	T_1	T_2	f_i	“FGM copula”	“Gumbel copula”
1	1	1	1	a	$pS_1S_2[1 + \varphi_D(1 - S_1)(1 - S_2)]$	$pS_1S_2Q_1$
2	1	1	0	b	$pS_1(1 - S_2)[1 - \varphi_D(1 - S_1)S_2]$	$pS_1[1 - S_2Q_1]$
3	1	0	1	c	$p(1 - S_1)S_2[1 - \varphi_D S_1(1 - S_2)]$	$pS_2[1 - S_1Q_1]$
4	1	0	0	[d]	$p(1 - S_1)(1 - S_2)[1 + \varphi_D S_1S_2]$	$p[1 - S_1 - S_2 + S_1S_2Q_1]$
5	0	1	1	e	$(1 - p)(1 - E_1)(1 - E_2)[1 + \varphi_{ND}E_1E_2]$	$(1 - p)(1 - E_1)(1 - E_2)Q_2$
6	0	1	0	f	$(1 - p)(1 - E_1)E_2[1 - \varphi_{ND}E_1(1 - E_2)]$	$(1 - p)(1 - E_1)[1 - (1 - E_2)Q_2]$
7	0	0	1	g	$(1 - p)E_1(1 - E_2)[1 - \varphi_{ND}E_2(1 - E_1)]$	$(1 - p)(1 - E_2)[1 - (1 - E_1)Q_2]$
8	0	0	0	[h]	$(1 - p)E_1E_2[1 + \varphi_{ND}(1 - E_1)(1 - E_2)]$	$(1 - p)[E_1 + E_2 - 1 + (1 - E_1)(1 - E_2)Q_2]$

Observe that; $Q_1 = \exp(-\phi_D \log S_1 \log S_2)$, $Q_2 = \exp(-\phi_{ND} \log(1 - E_1) \log(1 - E_2))$

3. Bayesian Approach

For a Bayesian analysis of the proposed models, we consider different Beta prior distributions on the prevalence, performance measure parameters (sensitivities and specificities) and the copula parameters. In some cases, we could have some prior information on the parameters from experts in diagnostic medical tests or from previous studies on the subject.

For a Bayesian analysis of the models, we assumed positive dependence between the diagnostic tests in the same way as it was considered by Dendukuri & Joseph (2001) (therefore $P(\varphi < 0) = 0$ and $P(\psi < 0) = 0$) we could assume uniform $U(a,b)$ as non-informative prior distributions and $Beta(a,b)$ distributions

for the informative situation for FGM and Gumbel dependence parameters and for prevalence and performance test parameters. If we need to elicit informative prior distributions for binary covariance, we could use the Generalized Beta(a,b) distribution in the same way that was considered by Martinez et al. (2005). For the non-informative case the Uniform U(0,1) distribution should be a good option.

Usually, we do not have any kind of information about the copula parameters, that is, for both copula dependence parameters. In this case, we used the procedure developed by Tovar (2012) to obtain the prior hyperparameters and we assume that the dependence takes values within of some interval (θ_1, θ_2) within of parametric space. In this way, if we assumed that the dependence is weak, the parameter could belong to the interval $(0, 1/4)$; when the dependence is moderate the parameter should be in to the interval $(1/4, 3/4)$ and when the dependence is strong, the parameter should be in to the interval $(3/4, 1)$. To obtain the hyperparameter values, we take the midpoint of the interval as the mean $E(\theta)$ and we apply the Chebychev's inequality to approximate the variance $V(\theta)$, as follows:

$$\begin{aligned} P(|\theta - E(\theta)| \geq k\sigma) &\leq \frac{1}{k^2} = \gamma \\ P([\theta - E(\theta)]^2 \geq k^2\sigma^2) &\leq \gamma \\ P(\alpha[\theta - \theta_0]^2 \geq \sigma^2) &\leq \gamma \end{aligned} \tag{8}$$

where γ is the prior probability of θ do not belong to the constructed interval.

Therefore, the variance will be a function of the prior established probability to interval values of the unknown quantity and the distance between θ_0 and a percentile of the distribution. If it is replaced θ by some of the known values θ_1 or θ_2 in the equation (8) it is easy to obtain a approximated value for the variance of the Beta prior distribution, as follows;

$$\sigma^2 \leq \alpha[\theta_1 - \theta_0]^2 \cong \frac{ab}{(a+b)^2(a+b+1)} \tag{9}$$

And as the mean $\theta_0 = E(\theta)$ and the variance $\sigma^2 = V(\theta)$ can be written as functions of the Beta prior hyperparameters, it is necessary to solve a system of two equations with two unknowns to find values of a and b i.e:

$$\begin{aligned} \omega &= \frac{\theta_0}{(1 - \theta_0)} \\ a &= \omega b \\ b &= \frac{\omega - [(\omega + 1)^2\sigma^2]}{(\omega^3 + 3\omega^2 + 3\omega + 1)\sigma^2} \end{aligned} \tag{10}$$

In this way, assuming $\gamma = 0.05$ in (8), for the FGM and Gumbel dependence parameters we have evaluated a Beta(17, 122) distribution, a Beta(39.5, 39.5) distribution and a Beta(122, 17) distribution as informative prior distributions and finally we have selected as selection criteria the Deviance Information Criteria DIC Spiegelhalter, Thomas, Best & Lunn (2003) obtained within the WinBUGS

environment and a heuristic procedure that assumes two criteria: quality in the convergence of the MCMC procedure and concentration of the posterior distribution using the coefficient of variation (CV). The best model should have the lower DIC, the best performance in MCMC convergence and highest concentration around the posterior mean (lowest CV).

We have seven parameters to be estimated, two sensitivities, two specificities, one prevalence, one dependence parameter for diseased individuals and another one for non-diseased individuals. If we assume a design with the presence of verification bias, we have only four degrees of freedom for the estimation process and if we assume a design without verification bias, we have six information components. Therefore, in both cases the model is non-identified. Using a classical approach, the problem has been addressed giving fixed values to a subset of parameters and estimating the remaining unconstrained parameters (Vacek 1985), but since all parameters are typically unknown, the division into constrained and unconstrained sets is often quite arbitrary. Since the Bayesian paradigm some authors as Joseph, Gyorkos & Coupal (1995), have proposed to construct informative prior distribution over a subset or over all unknown quantities. In accordance with Dendukuri & Joseph (2001), informative priors would be needed on at least as many parameters as would be constrained when using the most frequent approach. In this approach, the prior information is used to distinguish between the numerous possible solutions for the non-identifiable problem. This approach is approximately numerically equivalent to the most frequent approach when a degenerate (point mass) distribution is used that matches the constrained parameter values and diffuse prior distributions are used for the non-constrained parameters. In order to treat the non-identifiability problem, first, we assume informative prior distributions over the subset of dependence parameters and non-informative prior distributions on prevalence and performance test parameters and next, we assume informative prior distributions on all set of parameters in accordance with what was suggested by Joseph et al. (1995).

As the posterior distributions do not have closed forms, we have used MCMC methods, especially Metropolis-Hastings algorithm to obtain estimates for the parameters. For all models, 500,000 Gibbs samples were simulated from the conditional distributions. From these generated samples, we discarded the first 50,000 samples to eliminate the effect of the initial values and we also considered a spacing of 100. Convergence of the algorithm was verified graphically and also using standard existing methods implemented in the software CODA (Best, Cowles & Vines 1995).

4. Examples

4.1. Cancer Data

As a first example, we have used a data set introduced by (Smith et al. 1997). They screened 19,476 men for prostate cancer using Digital Rectal Examination (DRE) and Prostate-Specific Antigen (PSA) in serum. The PSA level was consid-

ered suspicious for cancer if it exceeded 4.0 ng/ml. Subjects with positive results on either DRE or PSA were submitted to an ultrasound guided needle biopsy test which was considered as “gold standard”. This data set obtained under verification bias is related to approximately 20,000 individuals, as such, it may be considered as a large sample size.

For prior distribution elicitation, we have used the results introduced by Böhning & Patilea (2008). We get the values for the δ and λ indexes and from these results, we estimated the quantities d and h of non-verified subjects given in Table 1. (See Table 4).

TABLE 4: Estimated values for the dependence indexes and quantities of non-verified individuals using Böhning’s results. The values in brackets were calculated using δ_i index, the another one using λ_i index.

	Diseased subjects $\lambda_1 = 2.42, \delta_1 = 3.08$			Non-diseased subjects $\lambda_0 = 2.40, \delta_0 = 3.03$		
	DRE+	DRE-	Total	DRE+	DRE-	Total
PSA+	189	292	481	141	755	896
PSA-	145	1431[691]	1576[836]	1002	15521[16261]	16523[17263]
Total	334	1723[983]	2057[1317]	1143	16276[17016]	17419[18159]

Using the data in Table 4 we assumed prior independence between the components of the parameter vector $[\theta_1 = S_1, \theta_2 = S_2, \theta_3 = E_1, \theta_4 = E_2, \theta_5 = p]$ to obtain estimates and intervals where it is possible assume to find each component with a probability $1 - \gamma = 0.95$. (See Table 5).

TABLE 5: Informative prior distribution hyperparameters for performance test parameters, prevalence and covariance (Martinez’s prior informative distributions for ψ).

PARAMETER	INTERVAL	E(θ)	a_θ	b_θ
S_1	0.236 - 0.365	0.3006	303	704
S_2	0.162 - 0.254	0.2080	324	1232
E_1	0.949 - 0.951	0.950	902500	47500
E_2	0.934 - 0.937	0.9355	501758	34595
p	0.068 - 0.106	0.0866	379	4002
ψ_D	0.004659 - 0.004719	0.004689	486303	103225102
ψ_{ND}	0.080 - 0.133	0.1722	289	2421

Assuming prior independence, for each interval we take the midpoint of each interval as the expected value of the prior distribution and we use the Chebychev inequality to get approximations for the variance of each parameter in the way that was described in Section 3 and we obtained the hyperparameter values that appear in Table 5. For this set of parameters we have used U(0,1) distributions as non-informative priors.

To elicit binary covariance prior distributions, we have used the results obtained by Martinez et al. (2005). They estimated the covariance parameter for the same cancer data under a Bayesian approach assuming non-informative prior distributions for ψ_D and ψ_{ND} . We have used the 95% credible regions obtained by them and we applied the same procedure employed with the test parameters

and prevalence. As non-informative distributions we have used GenBeta(1/2, 1/2) distributions.

For the copula parameters $\theta_2 = [\varphi_D, \varphi_{ND}, \phi_D, \phi_{ND}]$ we assumed the Beta distributions Beta(17, 122), Beta(39.5, 39.5) and Beta(122, 17) as prior distributions and Uniform U(0,1) as non-informative prior distributions. To address the lack identifiability problem of we have putting informative prior distributions on a subset or on the complete set of parameters considering two set of models as follows:

- Set 1 of models: informative prior distribution for the copula parameters and non-informative prior distributions for the prevalence and test parameters
- Set 2 of models: informative prior distributions for all parameters (See Table 6)

TABLE 6: Bayesian posterior summaries obtained by analyzing the data considering independence between tests assumption and different dependence structures. (Posterior means and 95% credible intervals (95% CrI) for each parameter of interest).

		Set 1 of models				Set 2 of models	
Model	Parameter	Means	95% CrI	Model	Parameter	Means	95% CrI
$M_{1,1}$ $DIC = 180.4$	S_1	0.567	0.529 - 0.605	$M_{2,1}$ $DIC = 337.1$	S_1	0.258	0.252 - 0.264
	S_2	0.394	0.363 - 0.394		S_2	0.226	0.208 - 0.244
	E_1	0.952	0.950 - 0.954		E_1	0.948	0.946 - 0.950
	E_2	0.946	0.943 - 0.949		E_2	0.947	0.944 - 0.950
	p	0.044	0.041 - 0.047		p	0.080	0.075 - 0.085
$M_{1,2}$ $DIC = 55.2$	ψ_D	0.0316	0.019 - 0.046	$M_{2,2}$ $DIC = 54.4$	ψ_D	0.046	0.037 - 0.055
	ψ_{ND}	0.005	0.004 - 0.006		ψ_{ND}	0.005	0.004 - 0.006
	S_1	0.470	0.380 - 0.548		S_1	0.295	0.274 - 0.316
	S_2	0.335	0.273 - 0.393		S_2	0.211	0.196 - 0.227
	E_1	0.951	0.948 - 0.955		E_1	0.950	0.950 - 0.950
	E_2	0.937	0.933 - 0.940		E_2	0.936	0.935 - 0.936
	p	0.051	0.044 - 0.062		p	0.082	0.076 - 0.088
$M_{1,3}$ $DIC = 156.5$	φ_D	0.156	0.136 - 0.176	$M_{2,3}$ $DIC = 225.5$	φ_D	0.135	0.123 - 0.148
	φ_{ND}	0.040	0.036 - 0.044		φ_{ND}	0.041	0.039 - 0.043
	S_1	0.538	0.480 - 0.595		S_1	0.320	0.300 - 0.343
	S_2	0.384	0.339 - 0.430		S_2	0.225	0.209 - 0.242
	E_1	0.952	0.948 - 0.955		E_1	0.950	0.949 - 0.950
	E_2	0.937	0.933 - 0.941		E_2	0.936	0.935 - 0.936
	p	0.045	0.040 - 0.050		p	0.074	0.069 - 0.079
$M_{1,4}$ $DIC = 192.7$	ϕ_D	0.120	0.072 - 0.179	$M_{2,4}$ $DIC = 294.4$	ϕ_D	0.047	0.028 - 0.072
	ϕ_{ND}	0.017	0.010 - 0.026		ϕ_{ND}	0.018	0.011 - 0.027
	S_1	0.593	0.540 - 0.645		S_1	0.330	0.307 - 0.353
	S_2	0.424	0.379 - 0.469		S_2	0.228	0.211 - 0.245
	E_1	0.952	0.948 - 0.955		E_1	0.950	0.949 - 0.951
	E_2	0.937	0.933 - 0.941		E_2	0.936	0.935 - 0.936
	p	0.040	0.037 - 0.044		p	0.072	0.0671 - 0.0771

$M_{j,1}$, $j = 1, 2$: Models under assumption of independence between tests
 $M_{j,2}$, $j = 1, 2$: Covariance parameters with informative prior distributions
 $M_{j,3}$, $j = 1, 2$: FGM dependence parameters with Beta(122, 17) prior distributions
 $M_{j,4}$, $j = 1, 2$: Gumbel dependence parameters with Beta (17, 122) prior distributions

From the results in Table 6, we observe that in this example with a large sample size (almost 20,000 individuals), we have great differences in the posterior summaries of interest, especially for the sensitivities S_ν , $\nu = 1, 2$ of the tests

considering different priors for the parameters and different modeling structures. It is also interesting to observe that the specificities E_ν $\nu = 1, 2$, that is, the probabilities of negative tests given that the individuals are not diseased, are almost not affected by the different priors and different modeling structures in presence or not of an dependence parameter. These results could be of great interest for medical diagnostic tests.

We also observe a large variability on the obtained DIC values considering each assumed model. The smallest DIC values are obtained for the class of models with a bivariate binary structure.

4.2. Urinary Tract Infection (UTI)

In this example, we consider a data set introduced by Ali et al. (2007) who evaluated a fast method to detect urinary tract infection. In this case, we can suspect an association between tests, since the results of the tests are more likely to be positive when the individual has a greater presence of infection. The authors considered the presence of nitrites ($N = test1$), and the levels of leukocyte esterase in urine ($LE = test2$) as screening tests and a bacterial culture as the confirmatory test. They applied the three methods in 132 children of both genders with ages ranging from three days to 11 years. The obtained performance test and prevalence estimates were compared with those obtained in other five studies. Since one of those studies had incomplete data, we only considered the results of the four complete studies to elicit our prior distributions. For each estimated parameter, we calculated the mean and variance of the results obtained in the five studies (including Ali's study) and used them as prior means and variances of the parameters. Thus, the informative prior distributions for prevalence and performance test parameters are given by:

$$S_1 \sim Beta(4.15, 4.5), \quad S_2 \sim Beta(15.7, 2.4)$$

$$E_1 \sim Beta(0.5, 13), \quad E_2 \sim Beta(8.3, 2.8)$$

and

$$p \sim Beta(2.1, 22.3)$$

For copula and covariance parameters, we assume the same informative priors used for copula parameters considered in the first example. We also assume uniform $U(0,1)$ prior distributions for the performance test parameters as non-informative priors and applied the same procedure for the estimation process used in the cancer data example. The results obtained are given in Table 7.

In this example with a small sample size, but not including missing data, we observe from Table 7, that the sensitivities S_ν $\nu = 1, 2$ were not greatly affected by the choice of prior distributions (informative or not) and modeling structures, but the specificities E_ν $\nu = 1, 2$ have a great variability considering the different modeling structures. We also observe that the prevalences have similar posterior summaries considering each model and the DIC values do not present great differences for each modeling structure.

TABLE 7: Bayesian posterior summaries obtained by analyzing the data considering independence between tests assumption and different dependence structures. (Posterior means and 95% credible intervals (95% CrI) for each parameter of interest).

		Set 1 of models				Set 2 of models		
Model	Parameter	Means	95% CrI	Model	Parameter	Means	95% CrI	
$M_{1,1}$ $DIC = 36.6$	S_1	0.387	0.318 - 0.457	$M_{2,1}$ $DIC = 40.3$	S_1	0.387	0.318 - 0.457	
	S_2	0.855	0.803 - 0.901		S_2	0.855	0.803 - 0.901	
	E_1	0.875	0.799 - 0.935		E_1	0.769	0.682 - 0.846	
	E_2	0.513	0.402 - 0.625		E_2	0.544	0.438 - 0.648	
	p	0.673	0.616 - 0.728		p	0.625	0.568 - 0.679	
$M_{1,1}$ $DIC = 36.4$	ψ_D	0.029	0.016 - 0.048	$M_{2,1}$ $DIC = 47.9$	ψ_D	0.028	0.015 - 0.049	
	ψ_{ND}	0.036	0.014 - 0.067		ψ_{ND}	0.077	0.049 - 0.110	
	S_1	0.384	0.287 - 0.484		S_1	0.392	0.298 - 0.489	
	S_2	0.847	0.767 - 0.912		S_2	0.857	0.785 - 0.916	
	E_1	0.870	0.756 - 0.949		E_1	0.702	0.582 - 0.809	
$M_{1,2}$ $DIC = 52.6$	E_2	0.567	0.424 - 0.702	$M_{2,2}$ $DIC = 40.0$	E_2	0.541	0.420 - 0.660	
	p	0.672	0.590 - 0.748		p	0.583	0.505 - 0.658	
	φ_D	0.050	0.027 - 0.794		$M_{2,2}$ $DIC = 40.0$	φ_D	0.053	0.028 - 0.085
	φ_{ND}	0.161	0.104 - 0.208			φ_{ND}	0.068	0.025 - 0.130
	S_1	0.392	0.289 - 0.487			S_1	0.385	0.289 - 0.487
S_2	0.855	0.783 - 0.915	S_2	0.845		0.764 - 0.911		
E_1	0.681	0.556 - 0.795	E_1	0.866		0.753 - 0.948		
$M_{1,3}$ $DIC = 42.6$	E_2	0.610	0.479 - 0.735	$M_{2,3}$ $DIC = 55.3$	E_2	0.578	0.433 - 0.716	
	p	0.583	0.505 - 0.658		p	0.672	0.590 - 0.748	
	ϕ_D	0.118	0.071 - 0.175		$M_{2,3}$ $DIC = 55.3$	ϕ_D	0.118	0.071 - 0.175
	ϕ_{ND}	0.119	0.072 - 0.177			ϕ_{ND}	0.121	0.073 - 0.179
	S_1	0.387	0.290 - 0.488			S_1	0.392	0.298 - 0.490
S_2	0.847	0.766 - 0.913	S_2	0.857		0.784 - 0.916		
E_1	0.864	0.751 - 0.946	E_1	0.679		0.553 - 0.792		
$M_{1,3}$ $DIC = 42.6$	E_2	0.576	0.431 - 0.715	$M_{2,3}$ $DIC = 55.3$	E_2	0.625	0.494 - 0.748	
	p	0.672	0.590 - 0.748		p	0.582	0.504 - 0.658	

$M_{j,1}, j = 1, 2$: Models under assumption of independence between tests
 $M_{j,2}, j = 1, 2$: Models using $GenBeta(39.5, 39.5)$ prior distributions for the covariance parameters
 $M_{j,3}, j = 1, 2$: Models taken $GenBeta(122, 17)$ prior distributions for the association FGM parameter
 $M_{j,4}, j = 1, 2$: Models with $Beta(17, 122)$ prior distributions for association Gumbel parameter

Considering DIC as discrimination criteria, we could assume a model with independence between the diagnostic tests considering informative or non-informative prior distributions or a model with dependence between tests given by a bivariate Bernoulli distribution (small DIC and similar performance test parameter estimates).

In this case, the copula parameter in non-diseased individuals presents an important change when we use informative priors over all parameters, while in the other group it remains unchanged. The specificity of the test N (test1) shows changes in the three models whether we use or do not use informative priors over the vector of non-dependence parameters. For the binary covariance and Gumbel models, the E_1 estimate with informative priors is lower than in the other models while in FGM model we observed an opposite behavior. When we have small sample size, the FGM model shows a more unstable behavior in the estimation of association parameter for non-diseased individuals. The DIC values for the different models do not show important changes. It is important, to observe that the DIC value of the FGM model with informative priors over all parameters is very similar with the DIC value of the Gumbel model when we use non-informative

priors over test parameters. On the other hand, the DIC value obtained assuming non-informative priors over test parameters in one model is similar with those obtained using informative priors over complete set in the other one. It is also interesting to see that the behavior of the FGM model with small sample size data is similar to the behavior observed in the Gumbel model when we have a large sample size.

5. Conclusion and Remarks

The main goal of this paper was to develop a Bayesian procedure to estimate the prevalence, performance test and copula parameters of two diagnostic tests in presence of verification bias and considering the dependence between test results.

We proposed the use of copula structure models to get the estimation of the parameters under dependence assumption and specifically, we have used the Farlie Gumbel Morgenstern (FGM) and the Gumbel copula models to compare the obtained results with a model under independence assumption between tests and another one assuming dependent binary tests in designs that consider two diagnostic tests with continuous outcome for screening, a perfect “gold standard” and verification bias presence. The estimation model obtained under verification bias presence, implies a lack of identifiability problem, because we have more parameters than informative pieces in the likelihood function. Given that, our approach considers the continuous dependence structure in the data but the estimation process is made with the binary observations in presence of verification bias, we consider that to estimate the parameters under the Bayesian approach is easier than under the frequentist approach, because many times it is possible that we do not have the continuous values, for instance, when the measured continuous traits are non-observable (they are latent variables).

We illustrated the procedure using two published data sets: one with a large sample size and another one with a small sample size of individuals. In both cases, the better fit for the data was obtained assuming binary associated tests and taking the covariance as a parameter. The FGM model showed better fit when compared to the Gumbel copula, regardless the sample size. With a large sample size, the FGM model presented DIC values lower when it was fitted assuming non-informative prior distributions on test parameters and the estimates are very close with those obtained using maximum likelihood method, reflecting the effect that has the observed data in the estimation process.

However, to use informative prior on all parameters allow us to obtain sensitivity estimates with shorter credibility regions which is very good if we consider that within the large sample used, the true positives are a small part. The previous conclusion is reinforced by the results observed with the data of the small sample size, which the informative prior on all parameters gave better fit. With the Gumbel model, we obtained similar results with large sample size, but the use of non-informative prior distributions on the test parameters gave better fit with small sample size. For binary covariance models the choice of prior distribution plays an important role in the estimation procedure, especially with large sample

sizes, where the posterior summaries of interest do not have important changes assuming informative or non-informative prior distributions. With small sample sizes and binary covariance structure, we observed better fit assuming non-informative prior distributions on the test parameters and informative prior distributions on covariance parameter.

It is important to point out that we could consider other copula families introduced in the literature to model dependence between diagnostic tests. A special case is given by the Clayton copula which is useful when the dependence is mainly concentrated in the lower tail or the Frank copula which is radial symmetric. The use of these other copulas in dependent diagnostic tests will be the goal of a future work, since an appropriate choice is essential in order to get an optimal result.

Acknowledgments

The second author was supported by grants of Conselho Nacional de Pesquisas (CNPq), Brazil. This paper is a part of the doctoral in statistics project of the first author.

[Recibido: noviembre de 2011 — Aceptado: abril de 2012]

References

- Ali, S., Moodambail, A., Hamrah, E., Bin-Nakhi, H. & Sadeq, S. (2007), 'Reliability of rapid dipstick test in detecting urinary tract infection in symptomatic children', *Kuwait Medical Journal* **39**, 36–38.
- Amblard, C. & Girard, S. (2002), 'Symmetry and dependence properties within a semiparametric family of bivariate copulas', *Journal of Non-parametric Statistics* **14**, 715–727.
- Amblard, C. & Girard, S. (2005), 'Estimation procedures for semiparametric family of bivariate copulas', *Journal of Computational and Graphical Statistics* **14**, 363–377.
- Amblard, C. & Girard, S. (2008), 'A new extension of bivariate FGM copulas', *Metrika* **70**, 1–17.
- Best, N., Cowles, M. & Vines, S. (1995), *CODA: Convergence diagnosis and output analysis software for Gibbs sampling output, version 0.3*; MRC Biostatistics Unit, Cambridge, U.K.
- Bohning, D. & Patilea, V. (2008), 'A capture-recapture approach for screening using two diagnostic tests with availability of disease status for the positives only', *Journal of the American Statistical Association* **103**, 212–221.
- Brenner, H. (1996), 'How independent are multiple diagnosis classifications?', *Statistics in Medicine* **15**, 1377–1386.

- Dendukuri, N. & Joseph, L. (2001), 'Bayesian approaches to modelling the conditional dependence between multiple diagnostic tests', *Biometrics* **57**, 158–167.
- Georgiadis, M., Johnson, W. & Gardner, I. (2003), 'Correlation adjusted estimation of sensitivity and specificity of two diagnostic tests', *Journal of the Royal Statistical Society: Series C (Applied Statistics)* **52**, 63–76.
- Gumbel, E. J. (1960), 'Bivariate exponential distributions', *Journal of the American Statistical Association* **55**, 698–707.
- Joseph, L., Gyorkos, T. W. & Coupal, L. (1995), 'Bayesian estimation of disease prevalence and the parameters of diagnostic tests in the absence of a gold standard', *American Journal of Epidemiology* **141**, 263–272.
- Martinez, E., Achcar, J. & Louzada, N. (2005), 'Bayesian estimation of diagnostic tests accuracy for semi-latent data with covariates', *Journal of Biopharmaceutical Statistics* **15**, 809–821.
- Nelsen, R. (1999), *An Introduction to Copulas*, Springer Verlag, New York.
- Qu, Y. & Hadgu, A. (1998), 'A model for evaluating sensitivity and specificity for correlated diagnostic test in efficacy studies with an imperfect reference test', *Journal of the American Statistical Association* **93**, 920–928.
- Smith, D., Bullock, A. & Catalona, W. (1997), 'Racial differences in operating characteristics of prostate cancer screening tests', *Journal of Urology* **158**, 1861–1865.
- Spiegelhalter, D., Thomas, A., Best, N. & Lunn, D. (2003), *Winbugs User Manual version 1.4*, MRC Biostatistics Unit, Cambridge, U.K.
- Thibodeau, L. (1981), 'Evaluating diagnostic tests', *Biometrics* **37**.
- Torrance-Rynard, V. & Walter, S. (1997), 'Effects of dependent errors in the assessment of diagnostic tests performance', *Statistics in Medicine* **16**.
- Tovar, J. R. (2012), 'Eliciting beta prior distributions for binomial sampling', *Revista Brasileira de Biometria* **30**, 159–172.
- Vacek, P. (1985), 'The effect of conditional dependence on the evaluation of diagnostic tests', *Biometrics* **41**.
- Walter, S. & Irwig, L. (1988), 'Estimation of test error rates disease prevalence and relative risk from misclassified data: a review', *Journal of Clinical Epidemiology* **41**.