Journal of Statistical Planning and Inference 🛚 (💵 🖿) 💵 – 💵



Contents lists available at ScienceDirect

Journal of Statistical Planning and Inference

journal homepage: www.elsevier.com/locate/jspi

Array-based schemes for group screening with test errors which incorporate a concentration effect

Y.G. Habtesllassie^a, Linda M. Haines^{b,*}, H.G. Mwambi^a, J.W. Odhiambo^c

^a School of Mathematics, Statistics and Computer Science, University of KwaZulu-Natal, Pietermaritzburg, South Africa ^b Department of Statistical Sciences, University of Cape Town, Private Bag X3, Rondebosch 7701, South Africa

^c Office of the Vice-Chancellor, Strathmore University, Nairobi, Kenya

ARTICLE INFO

Article history: Received 30 August 2014 Received in revised form 26 May 2015 Accepted 27 May 2015 Available online xxxx

Keywords: Array Group testing Factors False negatives False positives Sensitivity Specificity

ABSTRACT

Group screening is widely used as an efficient method for identifying samples or factors from a large population that are in some sense active. The focus in the present paper is on screening blood samples for infectious diseases when errors in testing are present. Specific attention is given to the introduction of a concentration effect, that is to settings in which the error in testing a group of blood samples depends on the number of samples in that group which are infected. Four array-based group screening schemes, the Dorfman, the AND, the OR and a modification of the AND scheme, are considered and their performance appraised by deriving explicit formulae for the expected number of tests, the expected number of false negatives and the expected number of false positives. The results are illustrated by means of two examples. As an aside, relationships complementary to those derived in the context of blood screening are developed within the area of group factor screening.

© 2015 Elsevier B.V. All rights reserved.

urnal of atistical planning

1. Introduction

In 1943 Dorfman introduced an efficient procedure for the screening of a large number of blood samples for a rare disease. Specifically, the scheme involves testing pooled blood samples for the disease and then, in a second stage, testing individual samples from those pools which tested positive for the disease in the first stage. In a complementary paper in 1961, Watson examined the related problem of group factor screening in cases of effect sparsity and proposed an approach based on that of Dorfman (1943) which incorporates appropriately designed two-stage experiments and which, additionally, accommodates the errors necessarily incurred in the test procedures. Since the publication of these two seminal papers, there has been considerable interest in issues relating to the broad area of group screening and to blood screening and group factor screening, case identification and prevalence estimation, and both of these have been extensively researched. In the context of case identification attention has focused, *inter alia*, on developing multi-stage and sequential schemes that are in some sense more effective than the Dorfman procedure (Kim et al., 2007) and on relaxing the somewhat stringent assumption that the probability that a blood sample is infected is constant (Bilder et al., 2010; McMahan et al., 2012).

There are surprisingly few reported studies on the effect of errors in testing in two-stage, and more generally in multistage, group screening procedures for the identification of samples which are in some way defective, such as infected blood

* Corresponding author. E-mail address: linda.haines@uct.ac.za (L.M. Haines).

http://dx.doi.org/10.1016/j.jspi.2015.05.009 0378-3758/© 2015 Elsevier B.V. All rights reserved.

2

ARTICLE IN PRESS

Y.G. Habtesllassie et al. / Journal of Statistical Planning and Inference 🛚 (IIIII) III-III

samples. A study based on ideas from the area of group factor screening is presented in the M.Sc. thesis of Habtesllassie (2004) and some related early work is noted in the review by Hughes-Oliver (2006). In a more recent paper, Kim et al. (2007) presented an in-depth investigation into the impact of test errors on hierarchical and two-dimensional square array schemes for case identification, in particular by evaluating the expected number of tests and selected error-based operating characteristics when the specificity and selectivity for group and individual tests are constant. Kim and Hudgens (2009) have extended this study to accommodate three-dimensional arrays. In addition, Hedt and Pagano (2008a,b) have devised an algorithm which accommodates errors in testing for square array pooling with halving at the retesting stage and have reported their findings in two working papers. In contrast, there has been a number of studies on the impact of errors in testing on the estimation of disease prevalence (Hughes-Oliver, 2006; Liu et al., 2012; Zhang et al., 2014). Interest has centered, *inter alia*, on the dilution effect, that is on settings in which false negatives which arise when a group comprising predominantly uninfected samples together with a very few infected samples tests negative. This effect was first discussed in the literature in the paper of Hwang (1976) and has been explored in more recent studies by, for example, Wein and Zenios (1996) and Hung and Swallow (1999). As noted earlier, errors in testing are necessarily embedded within the notion of group factor screening. Studies in this area are reviewed in Morris (2006) and have been restricted primarily to the Dorfman scheme and to extensions of that scheme to multi-stage and stepwise screening.

The aim of the present study is to extend the work reported in the literature on blood screening for the identification of infected samples in the presence of test errors by borrowing strength from the area of group factor screening and by taking cognizance of studies on prevalence estimation. More specifically, following Burns and Mauro (1987) and Habtesllassie (2004), the aim is to investigate blood screening schemes in which there is a concentration effect, that is one in which the errors in testing groups of blood samples depend on the number of infected individuals within the group. The paper is structured as follows. In Section 2 issues relating to the screening of blood samples are discussed. In particular four two-stage group screening schemes of interest are introduced and the expected number of tests, the expected number of false negatives and the expected number of false positives in the presence of test errors are evaluated explicitly for each scheme. The results are illustrated by means of examples in Section 3. In Section 4 the same issues as those discussed for blood screening are briefly revisited within the context of group factor screening and the differences which emerge, particularly in relation to the number of false negatives, are highlighted. Finally some broad conclusions and pointers for future research are given in Section 5. Note that, for ease of exposition, the setting in which blood samples are tested for an infectious disease is adopted here. However the discussion holds for a large number of group screening applications for case identification, including tests for defective items and in drug discovery. Note also that the terms concentration effect in the context of case identification and dilution effect in the context of prevalence estimation can be construed as being equivalent but that the emphasis in their interpretation is somewhat different. Thus, following Burns and Mauro (1987), the term concentration effect is used in the present study.

2. Blood screening

2.1. Preliminaries

Suppose that blood samples are arranged at random in a 2-dimensional array of cells. Four two-stage screening schemes based on this array are considered to be of interest here and are introduced as follows. The Dorfman scheme (Dorfman, 1943) involves the pooling and testing of samples in each row of the array and then the testing of individual samples in rows that test positive for the disease. The AND and the OR schemes, which were introduced in the paper by Phatarfod and Sudbury (1994) and the technical report by Langfeldt et al. (1997) respectively, are the same in the first stage and involve the pooling and testing of samples in each row and, independently, the pooling and testing of samples in each column. In the second stage of the AND scheme only individual samples lying at the intersection of rows and columns which tested positive in the first stage. Finally the scheme "square array without master pool testing" devised by Kim et al. (2007) is introduced and extended to a rectangular array. In this latter scheme all rows and columns are tested independently in the first stage. In the second stage, if at least one row and one column test positive, all cells at the intersection of positive rows are tested and similarly for columns. For all four schemes, if no rows or columns test positive in the first stage then the procedure stops.

The probability that an individual blood sample is infected, that is the prevalence of the disease, is taken to be a constant *p*, independent of the status of all other samples. In addition, this probability is assumed to remain unchanged on dilution, that is on pooling. Most importantly here, following notions for factor screening developed in the seminal paper of Watson (1961) and for group screening for defective items by Burns and Mauro (1987), errors in testing in stage one of each of the screening schemes of interest are taken to depend on the number of infected cells in the group or pool. A concentration effect is therefore introduced into the setting for blood screening. The errors in testing in the first stage can be quantified by introducing the generic probability

 $\pi_1^*(s, k) = P(a \text{ group of } k \text{ cells tests positive given that } s \text{ out of the } k \text{ cells are infected})$

where s = 0, ..., k. Thus the sensitivity of the test is taken to depend on the number of infected cells in the group and is given by $\pi_1^{\star}(s, k)$ for s = 1, 2, ..., k and the specificity of the test is given by $1 - \pi_1^{\star}(0, k)$. The errors in testing an individual

Y.G. Habtesllassie et al. / Journal of Statistical Planning and Inference II (IIIII) III-III

cell in the second stage are introduced through the probability

 $\pi_2^{\star}(a) = P(a \text{ cell with activity } a \text{ tests positive})$

where a = + denotes an infected cell and a = - an uninfected cell. The sensitivity of the test is thus given by $\pi_2^*(+)$ and the specificity by $\pi_2^N(-) = 1 - \pi_2^*(-)$. Finally, note that the errors in testing in the two stages for each of the blood screening schemes are assumed to be independent. The two stages themselves are however inextricably linked in that the number of cells tested in the second stage depends on the outcome of the first stage.

Three commonly used criteria within the context of errors in testing are examined in the present study in order to assess, compare and contrast the performances of the four screening algorithms (Langfeldt et al., 1997; Kim et al., 2007). Specifically, for a particular scheme *S* and a $k_1 \times k_2$ array of cells, these are the expected number of tests to be performed, $E_{T,S}(k_1, k_2)$, the expected number of false negatives, $E_{FN,S}(k_1, k_2)$, and the expected number of false positives, $E_{FP,S}(k_1, k_2)$. In the derivations which follow the Dorfman scheme is denoted by *D*, the AND scheme by *A*, the oR scheme by *O* and the scheme devised by Kim et al. (2007) and invoked here by *A*2. Note that the expected number of tests per cell is an operating characteristic usually referred to as the efficiency and that the expected number of false positives and the expected number of false negatives are termed the per-family error rate and the type II per-family error rate respectively (Kim et al., 2007). Note also that for a particular scheme other operating characteristics, such as the pooling positive predictive value (PPV), that is the probability that a cell is infected given that it tests negative, as delineated in Kim et al. (2007), can immediately be derived from the expected numbers of false negatives and false positives. Specifically, the pooling sensitivity and the pooling specificity for a scheme *S* are defined to be

$$S_e(S) = 1 - \frac{E_{FN,S}(k_1, k_2)}{k_1 k_2 p}$$
 and $S_p(S) = 1 - \frac{E_{FP,S}(k_1, k_2)}{k_1 k_2 (1 - p)}$

respectively and the pooling PPV and NPV are then given by

$$PPV(S) = \frac{pS_e(S)}{pS_e(S) + (1-p)(1-S_p(S))} \text{ and } NPV(S) = \frac{(1-p)S_p(S)}{(1-p)S_p(S) + p(1-S_e(S))}$$

2.2. Derivations

Consider first a generic derivation of the probabilities relating to the testing of a group G of k cells. Let G^* denote the event that the pooled group of cells tests positive. Then

$$P(G^{\star}) = \sum_{s=0}^{k} P(G^{\star}|s \text{ out of } k \text{ cells are infected}) P(s \text{ out of } k \text{ cells are infected})$$
$$= \sum_{s=0}^{k} \pi_{1}^{\star}(s, k) {\binom{k}{s}} p^{s} (1-p)^{k-s}.$$

Suppose also that *G* contains the cell *A* and let A^+ denote the event that the cell *A* is infected and A^- that it is not infected. Then

$$P(G^{\star}|A^{+}) = \sum_{s=0}^{k-1} P(G^{\star}|A^{+} \text{ and } s \text{ out of } k-1 \text{ cells are infected})P(s \text{ out of } k-1 \text{ cells are infected})$$
$$= \sum_{s=0}^{k-1} \pi_{1}^{\star}(s+1,k) {\binom{k-1}{s}} p^{s}(1-p)^{k-1-s}$$

and similarly

$$P(G^{\star}|A^{-}) = \sum_{s=0}^{k-1} \pi_{1}^{\star}(s,k) {\binom{k-1}{s}} p^{s} (1-p)^{k-1-s}.$$

Let R_i^*, C_j^* and A_{ij}^* denote the events that the *i*th row, the *j*th column and the (i, j)th cell respectively test positive and R_i^N and C_j^N denote the events that the *i*th row and the *j*th column respectively test negative. Further, let A_{ij}^+ and A_{ij}^- denote the events that the *i*th row and the *j*th column respectively test negative. Further, let A_{ij}^+ and A_{ij}^- denote the events that the blood sample in the (i, j)th cell A_{ij} is infected and uninfected respectively, with $P(A_{ij}^+) = p$ and $P(A_{ij}^-) = 1 - p \equiv q$. Probabilities associated with the events R_i^*, C_j^*, R_i^N and C_j^N , derived from the generic results above, are presented in Table 1 and provide the essential building blocks for the ensuing calculations. The notation introduced in Table 1 is used throughout the derivations which now follow.

4

ARTICLE IN PRESS

Y.G. Habtesllassie et al. / Journal of Statistical Planning and Inference 🛚 (💵 💵 – 💵

Table 1

Probabilities associated with the events R_i^{\star} , $R_i^{\rm N}$, C_i^{\star} and $C_i^{\rm N}$.

Probability	Notation	Expression	S _e and S _p
$P(R_i^{\star})$	π_R^*	$\sum_{s=0}^{k_2} \pi_1^{\star}(s, k_2) {\binom{k_2}{s}} p^s (1-p)^{k_2-s}$	$(1-S_p)q^{k_2}+S_e(1-q^{k_2})$
$P(R_i^{\rm N})$	π_R^{N}	$1 - \sum_{s=0}^{k_2} \pi_1^{\star}(s, k_2) {\binom{k_2}{s}} p^s (1-p)^{k_2-s}$	$S_p q^{k_2} + (1 - S_e)(1 - q^{k_2})$
$P(C_j^{\star})$	π_c^*	$\sum_{s=0}^{k_1} \pi_1^{\star}(s, k_1) {\binom{k_1}{s}} p^s (1-p)^{k_1-s}$	$(1-S_p)q^{k_1}+S_e(1-q^{k_1})$
$P(C_j^N)$	π_c^{N}	$1 - \sum_{s=0}^{k_1} \pi_1^{\star}(s, k_1) {\binom{k_1}{s}} p^s (1-p)^{k_1-s}$	$S_p q^{k_1} + (1 - S_e)(1 - q^{k_1})$
$P(R_i^\star A_{ij}^+)$	$\pi^*_{{ m extsf{R}} +}$	$\sum_{s=0}^{k_2-1} \pi_1^{\star}(s+1,k_2) \binom{k_2-1}{s} p^s (1-p)^{k_2-1-s}$	S _e
$P(R_i^{\rm N} A_{ij}^+)$	$\pi^{\mathrm{N}}_{\mathrm{R} +}$	$1 - \sum_{s=0}^{k_2-1} \pi_1^{\star}(s+1,k_2) \binom{k_2-1}{s} p^s (1-p)^{k_2-1-s}$	$1 - S_e$
$P(R_i^\star A_{ij}^-)$	$\pi^*_{R -}$	$\sum_{s=0}^{k_2-1} \pi_1^{\star}(s, k_2) \binom{k_2 - 1}{s} p^s (1-p)^{k_2 - 1 - s}$	$(1-S_p)q^{k_2-1}+S_e(1-q^{k_2-1})$
$P(R_i^{\rm N} A_{ij}^-)$	$\pi^{\mathrm{N}}_{R -}$	$1 - \sum_{s=0}^{k_2-1} \pi_1^{\star}(s, k_2) \binom{k_2 - 1}{s} p^s (1-p)^{k_2-1-s}$	$S_p q^{k_2 - 1} + (1 - S_e)(1 - q^{k_2 - 1})$
$P(C_j^{\star} A_{ij}^+)$	$\pi^*_{\mathcal{C} +}$	$\sum_{s=0}^{k_1-1} \pi_1^{\star}(s+1,k_1) \binom{k_1-1}{s} p^s (1-p)^{k_1-1-s}$	Se
$P(C_j^{\rm N} A_{ij}^+)$	$\pi^{\mathrm{N}}_{C +}$	$1 - \sum_{s=0}^{k_1-1} \pi_1^{\star}(s+1, k_1) \binom{k_1 - 1}{s} p^s (1-p)^{k_1-1-s}$	$1-S_e$
$P(C_j^\star A_{ij}^-)$	$\pi^*_{C -}$	$\sum_{s=0}^{k_1-1} \pi_1^{\star}(s,k_1) \binom{k_1-1}{s} p^s (1-p)^{k_1-1-s}$	$(1-S_p)q^{k_1-1}+S_e(1-q^{k_1-1})$
$P(C_j^{\rm N} A_{ij}^-)$	$\pi^{\mathrm{N}}_{C -}$	$1 - \sum_{s=0}^{k_1-1} \pi_1^{\star}(s, k_1) \binom{k_1 - 1}{s} p^s (1-p)^{k_1-1-s}$	$S_p q^{k_1 - 1} + (1 - S_e)(1 - q^{k_1 - 1})$

2.2.1. Dorfman

Consider first the Dorfman scheme for screening blood samples. Expressions for the criteria of interest can be derived by arguments similar to those used by Watson (1961) in the context of factor screening and more specifically by Langfeldt et al. (1997) in the context of blood screening with blockers. Thus the expected number of tests follows immediately as

 $E_{T,D}(k_1, k_2) = k_1 + k_1 k_2 P(R_i^*)$ = $k_1 + k_1 k_2 \pi_R^*$.

The expected number of true positives is derived by conditioning on the status of a particular cell and is given by

$$E_{TP,D}(k_1, k_2) = k_1 k_2 P(A_{ij}^+ \cap R_i^* \cap A_{ij}^*)$$

= $k_1 k_2 P(R_i^* | A_{ij}^+) P(A_{ij}^* | A_{ij}^+) P(A_{ij}^+)$
= $k_1 k_2 p \pi_{R_{l+}}^* \pi_2^*(+)$

and the expected number of false negatives is therefore

$$E_{FN,D}(k_1, k_2) = k_1 k_2 P(A_{ij}^+) - E_{TP,D}(k_1, k_2)$$

= $k_1 k_2 p \{1 - \pi_{R|+}^* \pi_2^*(+)\}.$

The expected number of false positives is similarly derived as

$$\begin{split} E_{FP,D}(k_1, k_2) &= k_1 k_2 \, P(A_{ij}^- \cap R_i^\star \cap A_{ij}^\star) \\ &= k_1 k_2 \, P(R_i^\star | A_{ij}^-) P(A_{ij}^\star | A_{ij}^-) \, P(A_{ij}^-) \\ &= k_1 k_2 \, (1-p) \, \pi_{R|-}^\star \, \pi_2^\star(-). \end{split}$$

These results are in accord with the results derived by Burns and Mauro (1987) for the special case of multiple testing for defective items corresponding to the Dorfman scheme.

2.2.2. The AND algorithm

The expected number of tests, the expected number of false negatives and the expected number of false positives for the AND and the OR schemes can be derived by invoking arguments similar to those for the Dorfman scheme. Specifically, the required expressions for the AND scheme can be obtained by replacing the event R_i^* in the derivations of the previous section with the event $R_i^* \cap C_i^*$ and by observing that the testing of rows and columns is independent so that $P(R_i^* \cap C_i^* | A_{ii}^+) =$

Y.G. Habtesllassie et al. / Journal of Statistical Planning and Inference I (IIII) III-III

 $P(R_i^{\star}|A_{ii}^+) P(C_i^{\star}|A_{ii}^+)$ and $P(R_i^{\star} \cap C_i^{\star}|A_{ii}^-) = P(R_i^{\star}|A_{ii}^-) P(C_i^{\star}|A_{ii}^-)$. Thus the expected number of tests is given by

$$\begin{split} E_{T,A}(k_1, k_2) &= k_1 + k_2 + k_1 k_2 P(R_i^{\star} \cap C_j^{\star}) \\ &= k_1 + k_2 + k_1 k_2 \left\{ P(R_i^{\star} | A_{ij}^+) P(C_j^{\star} | A_{ij}^+) P(A_{ij}^+) + P(R_i^{\star} | A_{ij}^-) P(C_j^{\star} | A_{ij}^-) P(A_{ij}^-) \right\} \\ &= k_1 + k_2 + k_1 k_2 \left\{ \pi_{R|+}^* \pi_{C|+}^* p + \pi_{R|-}^* \pi_{C|-}^* (1-p) \right\}, \end{split}$$

the expected number of false negatives by

$$\begin{split} E_{FN,A}(k_1, k_2) &= k_1 k_2 \{ P(A_{ij}^+) - P(A_{ij}^+ \cap (R_i^* \cap C_j^*) \cap A_{ij}^*) \} \\ &= k_1 k_2 \left\{ P(A_{ij}^+) - P(R_i^* | A_{ij}^+) P(C_j^* | A_{ij}^+) P(A_{ij}^* | A_{ij}^+) P(A_{ij}^+) \right\} \\ &= k_1 k_2 p \left\{ 1 - \pi_{R|+}^* \pi_{C|+}^* \pi_2^*(+) \right\} \end{split}$$

and the expected number of false positives by

$$\begin{split} E_{FP,A}(k_1, k_2) &= k_1 k_2 P(A_{ij}^- \cap (R_i^\star \cap C_j^\star) \cap A_{ij}^*) \\ &= k_1 k_2 P(R_i^\star | A_{ij}^-) P(C_j^\star | A_{ij}^-) P(A_{ij}^* | A_{ij}^-) P(A_{ij}^-) \\ &= k_1 k_2 (1-p) \pi_{R|-}^* \pi_{C|-}^* \pi_2^\star (-). \end{split}$$

2.2.3. The OR algorithm

Expressions for the criteria of interest in the OR scheme follow by replacing the event R_i^* in the derivations for the Dorfman scheme with the event $R_i^* \cup C_j^*$, by invoking the standard probability relationship $P(R_i^* \cup C_j^*) = P(R_i^*) + P(C_j^*) - P(R_i^* \cap C_j^*)$ and by using the results for the AND scheme. The results for the OR scheme can then be summarized succinctly as follows. Thus the expected number of tests is given by

$$\begin{split} E_{T,O}(k_1, k_2) &= k_1 + k_2 + k_1 k_2 P(R_i^{\star} \cup C_j^{\star}) \\ &= k_1 + k_2 + k_1 k_2 \{ P(R_i^{\star}) + P(C_j^{\star}) - P(R_i^{\star} \cap C_j^{\star}) \} \\ &= k_1 + k_2 + k_1 k_2 \{ \pi_R^{\star} + \pi_C^{\star} - [\pi_{R|+}^{\star} \pi_{C|+}^{\star} p + \pi_{R|-}^{\star} \pi_{C|-}^{\star} (1-p)] \}, \end{split}$$

the expected number of false negatives by

$$\begin{split} E_{FN,0}(k_1, k_2) &= k_1 k_2 \{ P(A_{ij}^+) - P(A_{ij}^+ \cap (R_i^* \cup C_j^*) \cap A_{ij}^*) \} \\ &= k_1 k_2 \left\{ P(A_{ij}^+) - P(R_i^* \cup C_j^* | A_{ij}^+) P(A_{ij}^* | A_{ij}^+) P(A_{ij}^*) \right\} \\ &= k_1 k_2 \left\{ P(A_{ij}^+) - \left[P(R_i^* | A_{ij}^+) + P(C_j^* | A_{ij}^+) - P(R_i^* \cap C_j^* | A_{ij}^+) \right] P(A_{ij}^* | A_{ij}^+) P(A_{ij}^*) \right\} \\ &= k_1 k_2 p \left\{ 1 - (\pi_{R|+}^* + \pi_{C|+}^* - \pi_{R|+}^* \pi_{C|+}^*) \pi_2^*(+) \right\} \end{split}$$

and the expected number of false positives by

$$\begin{split} E_{FP,O}(k_1, k_2) &= k_1 k_2 \, P(A_{ij}^- \cap (R_i^\star \cup C_j^\star) \cap A_{ij}^\star) \\ &= k_1 k_2 P(R_i^\star \cup C_j^\star | A_{ij}^-) \, P(A_{ij}^\star | A_{ij}^-) \, P(A_{ij}^-) \\ &= k_1 k_2 [P(R_i^\star | A_{ij}^-) + P(C_j^\star | A_{ij}^-) - P(R_i^\star \cap C_j^\star | A_{ij}^-)] \, P(A_{ij}^\star | A_{ij}^-) \, P(A_{ij}^-) \\ &= k_1 k_2 \, (1-p) \, (\pi_{R_{I-}}^\star + \pi_{C_{I-}}^\star - \pi_{R_{I-}}^\star \pi_{C_{I-}}^\star) \, \pi_2^\star(-). \end{split}$$

2.2.4. The A2 (k_1, k_2) algorithm

In evaluating the operating characteristics for the $A2(k_1, k_2)$ scheme, the event $R_i^* \cap C_j^*$, that is the AND scheme, the event $R_i^* \cap_{j=1}^k C_j^N$, that is the case in which the *i*th row of the array tests positive and all columns test negative and the event $\bigcap_{i=1}^{k_1} R_i^N \cap C_j^*$, that is the case in which the *j*th column tests positive and all rows test negative for $i = 1, ..., k_1$ and $j = 1, ..., k_2$, must necessarily be considered (Kim et al., 2007). The expected number of tests, the expected number of false negatives and the expected number of false positives associated with the event $R_i^* \cap C_j^*$ have already been derived in Section 2.2.2. It is therefore only necessary to consider the event $R_i^* \cap_{j=1}^{k_2} C_j^N$ with $i = 1, ..., k_1$. The evaluations associated with the event $\bigcap_{i=1}^k R_i^N \cap C_i^*$ with $i = 1, ..., k_1$.

with the event $\bigcap_{i=1}^{k_1} R_i^N \cap C_j^*$ with $j = 1, ..., k_2$, mirror those for the event $R_i^* \cap_{j=1}^{k_2} C_j^N$ and can then be written down. The requisite derivations for the A2 scheme are in fact somewhat intricate and, for conciseness, are presented in the Appendix. Thus, consider the term

- 5

$$\bar{S}_{(e,R)} = P(R_i^{\star} | \text{at least one } A_{ij}^+) = \frac{\sum_{s=1}^{k_2} \pi_1^{\star}(s, k_2) {\binom{k_2}{s}} p^s (1-p)^{k_2-1}}{(1-q^{k_2})}$$

Y.G. Habtesllassie et al. / Journal of Statistical Planning and Inference & (****) ***-***

which can be interpreted as the average sensitivity over rows for the stage one test procedure, together with the term $\bar{S}_{(e,C)}$ which is equal to $\bar{S}_{(e,R)}$ but with k_2 replaced by k_1 . It then follows that the expected number of tests is given by

$$\begin{split} E_{T,A2}(k_1,k_2) &= k_1 + k_2 + k_1 k_2 \left[P(R_i^{\star} \cap C_j^{\star}) + P\left(R_i^{\star} \bigcap_{j=1}^{k_2} C_j^{\mathsf{N}}\right) + P\left(\bigcap_{i=1}^{k_1} R_i^{\mathsf{N}} \cap C_j^{\star}\right) \right] \\ &= k_1 + k_2 + k_1 k_2 \left\{ \pi_{R|+}^{\star} \pi_{C|+}^{\star} p + \pi_{R|-}^{\star} \pi_{C|-}^{\star} (1-p) \right. \\ &+ \left[\pi_1^{\star}(0,k_2) - \bar{S}_{(e,R)} \right] (\pi_{C|-}^{\mathsf{N}} q)^{k_2} + \bar{S}_{(e,R)} (\pi_C^{\mathsf{N}})^{k_2} \\ &+ \left[\pi_1^{\star}(0,k_1) - \bar{S}_{(e,C)} \right] (\pi_{R|-}^{\mathsf{N}} q)^{k_1} + \bar{S}_{(e,C)} (\pi_R^{\mathsf{N}})^{k_1} \right\}, \end{split}$$

the expected number of false negatives by

$$E_{FN,A2}(k_1, k_2) = k_1 k_2 p \left\{ 1 - \left[P(R_i^{\star} \cap C_j^{\star} \cap A_{ij}^{\star} | A_{ij}^{+}) + P\left(R_i^{\star} \bigcap_{j=1}^{k_2} C_j^{\mathsf{N}} \cap A_{ij}^{\star} | A_{ij}^{+}\right) + P\left(\bigcap_{i=1}^{k_1} R_i^{\mathsf{N}} \cap C_j^{\star} \cap A_{ij}^{\star} | A_{ij}^{+}\right) \right] \right\}$$

= $k_1 k_2 p \left\{ 1 - \left[\pi_{R|+}^{\star} \pi_{C|+}^{\star} + \pi_{R|+}^{\star} (\pi_C^{\mathsf{N}})^{k_2 - 1} \pi_{C|+}^{\mathsf{N}} + \pi_{C|+}^{\star} (\pi_R^{\mathsf{N}})^{k_1 - 1} \pi_{R|+}^{\mathsf{N}}] \pi_2^{\star}(+) \right\}$

and the expected number of false positives by

$$\begin{split} E_{FP,A2}(k_1,k_2) &= k_1 k_2 \left[P(R_i^{\star} \cap C_j^{\star} | A_{ij}^{-}) + P\left(R_i^{\star} \bigcap_{j=1}^{k_2} C_j^{N} | A_{ij}^{-}\right) + P\left(\bigcap_{i=1}^{k_1} R_i^{N} \cap C_j^{\star} | A_{ij}^{-}\right) \right] P(A_{ij}^{\star} | A_{ij}^{-}) P(A_{ij}^{-}) \\ &= k_1 k_2 (1-p) \left\{ \pi_{R|-}^{\star} \pi_{C|-}^{\star} + (\pi_1^{\star}(0,k_2) - \bar{S}_{(e,R|-)})(\pi_{C|-}^{N} q)^{k_2 - 1} \pi_{C|-}^{N} + \bar{S}_{(e,R|-)}(\pi_{C}^{N})^{k_2 - 1} \pi_{C|-}^{N} \right. \\ &+ (\pi_1^{\star}(0,k_1) - \bar{S}_{(e,C|-)})(\pi_{R|-}^{N} q)^{k_1 - 1} \pi_{R|-}^{N} + \bar{S}_{(e,C|-)}(\pi_{R}^{N})^{k_1 - 1} \pi_{R|-}^{N} \right\} \pi_2^{\star}(-). \end{split}$$

2.2.5. Sensitivity and specificity: No dependence on the numbers of infected cells

Suppose now that the sensitivity and specificity associated with the rows and columns of a $k_1 \times k_2$ array do not depend on *s*, the number of infected cells. Thus let

$$\pi_1^{\star}(s, k) = S_e$$
 for $s = 1, 2, ..., k$ and $\pi_1^{\star}(0, k) = 1 - S_p$

for $k = k_1$ and k_2 . Then the probabilities and conditional probabilities associated with the events R_i^* and C_j^* follow immediately and are given in the last column of Table 1. In addition suppose that the sensitivity and specificity are the same for testing an individual cell as for testing the rows and columns, that is $\pi_2^*(+) = S_e$ and $\pi_2^*(-) = 1 - S_p$. For conciseness, following Kim et al. (2007), let

$$t(k) = (1 - S_p)q^k + S_e(1 - q^k).$$

Then, for the Dorfman scheme,

$$E_{T,D}(k_1, k_2) = k_1 + k_1 k_2 t(k_2)$$

$$E_{FN,D}(k_1, k_2) = k_1 k_2 p \{1 - S_e^2\}$$

$$E_{FP,D}(k_1, k_2) = k_1 k_2 (1 - p) (1 - S_p) t(k_2 - p)$$

in accord with the results of Kim et al. (2007), for the AND scheme,

$$E_{T,A}(k_1, k_2) = k_1 + k_2 + k_1 k_2 [pS_e^2 + (1-p) t(k_1 - 1) t(k_2 - 1)]$$

$$E_{FN,A}(k_1, k_2) = k_1 k_2 p \{1 - S_e^3\}$$

$$E_{FP,A}(k_1, k_2) = k_1 k_2 (1-p) (1 - S_p) t(k_1 - 1) t(k_2 - 1)$$
or the OB scheme

and for the OR scheme,

$$E_{T,O}(k_1, k_2) = k_1 + k_2 + k_1 k_2 \{t(k_1) + t(k_2) - [pS_e^2 + (1-p) t(k_1-1) t(k_2-1)]\}$$

$$E_{FN,O}(k_1, k_2) = k_1 k_2 p \left[1 - S_e^2 (2 - S_e)\right]$$

$$E_{FP,O}(k_1, k_2) = k_1 k_2 (1-p) (1 - S_p) [t(k_1-1) + t(k_2-1) - t(k_1-1)t(k_2-1)].$$

1)

Derivations for the $A2(k_1, k_2)$ scheme are again a little more involved than those for the Dorfman, the AND and the OR schemes and are detailed in the Appendix. Thus the expected number of tests, the expected number of false negatives and the expected number of false positives for the $A2(k_1, k_2)$ scheme are given by

$$\begin{split} E_{T,A2}(k_1,k_2) &= k_1 + k_2 + k_1 k_2 \left\{ S_e^2 p + t(k_1 - 1)t(k_2 - 1) (1 - p) \right. \\ &+ (1 - S_e - S_p) \left\{ \left[(1 - t(k_1 - 1)) q \right]^{k_2} + \left[(1 - t(k_2 - 1)) q \right]^{k_1} \right\} \\ &+ S_e \left\{ \left[(1 - t(k_1 - 1)) \right]^{k_2} + S_e \left[(1 - t(k_2 - 1)) \right]^{k_1} \right\} \right\}, \end{split}$$

$$\begin{split} E_{FN,A2}(k_1,k_2) &= k_1 k_2 p \left\{ 1 - S_e^3 - S_e^2 (1 - S_e) \left[(1 - t(k_1))^{k_2 - 1} + (1 - t(k_2))^{k_1 - 1} \right] \right\} \end{split}$$

Y.G. Habtesllassie et al. / Journal of Statistical Planning and Inference I (IIII) III-III

Table 2

Ordering of the group screening algorithms for expected numbers.

Number	Order
$ \begin{split} & E_{T,S}(k_1,k_2) \\ & E_{FN,S}(k_1,k_2) \\ & E_{FP,S}(k_1,k_2) \end{split} $	$0 > A2 \ge A$ $A > D > O \text{ and } A \ge A2$ $A < D < O \text{ and } A \le A2$

and

$$E_{FP,A2}(k_1, k_2) = k_1 k_2 (1-p)(1-S_p) \{t(k_1-1)t(k_2-1) + (1-S_e-S_p) \{[1-t(k_1-1)]^{k_2} q^{k_2-1} + [1-t(k_2-1)]^{k_1} q^{k_1-1} \}$$

+ $S_e \{[1-t(k_1)]^{k_2} + [1-t(k_2)]^{k_1} \} \}$

respectively.

Finally note that for perfect testing $S_e = S_p = 1$. Thus

$$\begin{split} E_{T,D}(k_1, k_2) &= k_1 + k_1 k_2 \; (1 - q^{k_2}) \\ E_{T,A}(k_1, k_2) &= k_1 + k_2 + k_1 k_2 \{1 - q^{k_1} - q^{k_2} + q^{k_1 + k_2 - 1}\} \\ E_{T,O}(k_1, k_2) &= k_1 + k_2 + k_1 k_2 \; \{1 - q^{k_1 + k_2 - 1}\} \\ E_{T,A2}(k_1, k_2) &= E_{T,A}(k_1, k_2) \end{split}$$

and $E_{FN,S}(k_1, k_2) = E_{FP,S}(k_1, k_2) = 0$ for all four schemes, that is for S = D, A, O and A2. These results are in accord with those presented, *inter alia*, in the papers by Langfeldt et al. (1997) and Kim et al. (2007), in the M.Sc. thesis of Habtesllassie (2004) and in the paper by Morris (1987) in the context of group factor screening.

2.2.6. Comparisons

Comparisons with respect to the performance of the four two-stage blood testing schemes of interest can now be made. Specifically, it is clear from the intrinsic nature of the associated algorithms that the expected number of tests for the or scheme is greater than that for the A2 scheme which in turn is greater than that for the AND scheme, that is $O > A2 \ge A$. The relation of the expected number of tests for the Dorfman scheme to this ordering is not however clear. It can also be deduced from the formulae presented in Sections 2.2.1–2.2.4 that the expected number of false negatives is greater for the AND scheme than that for the Dorfman scheme which in turn is greater than that for the OR scheme, that is A > D > O, and that the expected number of false negatives for the A2 scheme is less than that for the AND scheme, that is A > D > O, and that the expected number of false negatives, that is A < D < O and $A2 \ge A$. These observations are summarized in Table 2. It is now immediately clear from these relationships that if the expected numbers of tests, of false negatives and of false positives are all of interest then no single scheme yields optimal results. Rather, following Langfeldt et al. (1997), a costing of these operating characteristics is required and must necessarily be provided by the clinical researcher.

3. Examples

In order to construct examples it is first necessary to specify the explicit form of the probability that a group of *k* cells tests positive given that *s* cells are infected, that is $\pi_1^*(s, k)$, together with the values of the sensitivity and specificity of the individual tests, $\pi_2^*(+)$ and $\pi_2^*(-)$ respectively. Two forms for the probability $\pi_1^*(s, k)$, one taken from the case identification literature and one from the prevalence estimation literature, are introduced in the examples which follow and their attendant operating characteristics are explored.

3.1. Example 1

The form of $\pi_1^*(s, k)$ proposed by Burns and Mauro (1987) within the context of case identification is adopted in this example and is given by

$$\pi_1^*(s,k) = \alpha_1 + (1 - \alpha_1 - \alpha_2) \left(\frac{s}{k}\right)^{\gamma} \quad \text{for } s = 0, \dots, k,$$
(1)

where $0 \le \alpha_1, \alpha_2 \le 1, 0 < \alpha_1 + \alpha_2 < 1$ and $0 \le \gamma \le 1$. Note that the specificity for group testing is given by $\pi_1^*(0, k) = \pi_1^*(0) = \alpha_1$ and that $\pi_1^*(k, k) = 1 - \alpha_2$, thereby enabling formula (1) to be meaningfully calibrated. Note also that for $\gamma = 1, \pi_1^*(s, k)$ is linear in *s* and that as γ approaches 0, so $\pi_1^*(s, k)$ for $s = 1, 2, \ldots, k$ approaches the constant $\pi_1^*(k, k)$. Plots of the probability $\pi_1^*(s, k)$ given in formula (1) against *s* for $s = 0, \ldots, k, k = 8, \alpha_1 = 0.01, \alpha_2 = 0.1$ and selected values of γ are presented in Fig. 1(a) (see also Burns and Mauro (1987)).

ANNULL IN FILLU

Y.G. Habtesllassie et al. / Journal of Statistical Planning and Inference 🛚 (💵 🖿)



Fig. 1. Plots of $\pi_1^*(s, k)$ against *s* for k = 8, s = 0, 1, ..., 8, $\alpha_1 = 0.01$, $\alpha_2 = 0.1$ and (a) the Burns and Mauro (1987) sub-model with $\gamma = 0, 0.1, 0.25, 0.5, 0.75, 1.0$ and (b) the Hung and Swallow (1999) sub-model with $f_1 = 1.25$ and $f_2 = 0, 0.1, 1, 10, 100$.

Consider now evaluating the expected number of tests, the expected number of false negatives and the expected number of false positives for an 8 × 12 rectangular array over the range of prevalences $0 . Suppose that in order to model a concentration effect, the form of <math>\pi_1^*(s, k)$ given in (1) is adopted with $\alpha_1 = 0.01$, $\alpha_2 = 0.1$ and with k = 8 for testing the rows and k = 12 for testing the columns. Suppose also that the sensitivity and specificity of individual tests are taken to be $\pi_2^*(+) = 0.9$ and $\pi_2^N(-) = 0.99$ respectively and that $\gamma = 0.1$. Then, for these settings, plots of $E_{T,S}(8, 12)$, $E_{FN,S}(8, 12)$ and $E_{FP,S}(8, 12)$ against p for S = D, A, O and A2 and O are presented in Fig. 2 and plots of the pooling positive and pooling negative predictive values against <math>p for 0 and <math>0 respectively in Fig. 3.

Key features relating to the present example can be identified from the plots in Figs. 2 and 3, together with some additional minor calculations. Specifically, it is clear from the plots in Fig. 2 that the relationships between the four schemes, the Dorfman, the AND, the OR and the A2, for the expected number of tests, the expected number of false negatives and the expected number of false positives as established in Section 2.2.6 hold. In addition it is clear that orderings that could

ARTICLE IN PRESS

Y.G. Habtesllassie et al. / Journal of Statistical Planning and Inference I (IIIII) III-III



Fig. 2. Plots of (a) the expected number of tests, (b) the expected number of false negatives and (c) the expected number of false positives against *p* for $0 and the Burns and Mauro (1987) sub-model with <math>k_1 = 8$, $k_2 = 12$, $\alpha_1 = 0.01$, $\alpha_2 = 0.1$, $\pi_2^{\Lambda}(+) = 0.9$, $\pi_2^{\Lambda}(-) = 0.99$ and $\gamma = 0.1$.

not be quantified earlier change with p for small values of p, as indicated in Fig. 2. Note that no further changes in the ordering occurred from the end-point prevalence p = 0.05 in Fig. 2 up to p = 0.1. Note also that the expected number of tests, the expected number of false negatives and the expected number of false positives for the A2 scheme approach the

Y.G. Habtesllassie et al. / Journal of Statistical Planning and Inference 🛚 (💵 💷 – 💵



Fig. 3. Plots of (a) the pooling positive predicted values against *p* for 0 and (b) the pooling positive predicted values against*p* $for <math>0 and the Burns and Mauro (1987) sub-model with <math>k_1 = 8$, $k_2 = 12$, $\alpha_1 = 0.01$, $\alpha_2 = 0.1$, $\pi_2^{N}(-) = 0.9$, $\pi_2^{N}(-) = 0.99$ and $\gamma = 0.1$.

corresponding values for the AND scheme as *p* increases. In fact for prevalences, *p*, greater than 0.1 there is little advantage to be gained by invoking the A2 as opposed to the AND scheme. Finally note, more generally, that the differences between the expected numbers of tests, of false negatives and of false positives across certain schemes and prevalences are relatively small and that this could influence the choice of a suitable scheme.

The plots of the pooling PPV and the pooling NPV against p shown in Fig. 3 are particularly interesting since these values are of critical importance in testing for a given disease. Specifically, the pooling PPVs for the AND scheme are very close to 1 for all p in the range 0 , as compared with those for the Dorfman, the OR and the A2 schemes, thereby reflecting the fact that the expected numbers of false positives for the AND scheme are close to 0 and hence that the associated pooling specificities are close to 1. In contrast, for all the four schemes of interest the pooling NPVs are close to 1, ranging from 1 to 0.948 for prevalences <math>p from 0 to 0.1. Thus, on balance, it would seem that in such cases the AND scheme is to be preferred within the context of pooling PPVs and NPVs.

Clearly the above results relate to the present example. Examples with $\pi_1^*(s, k)$ given by formula (1) with the parameter γ ranging from 0 to 1 and with square and rectangular arrays were also investigated and were found to exhibit similar

Y.G. Habtesllassie et al. / Journal of Statistical Planning and Inference II (IIIII) III-III

overall features to those presented here but with different cut-off values and different expected numbers of tests, false negatives and false positives. For example, for an 80×80 array, the plots of the expected number of tests, false negatives and false positives against the prevalence for the AND and the A2 schemes are effectively the same, down to a prevalence of approximately 0.0015.

3.2. Example 2

Hung and Swallow (1999) introduced a sub-model for $\pi_1^*(s, k)$ of the form

$$\pi_{1}^{\star}(s,k) = \begin{cases} \alpha_{1} + (1-\alpha_{1}-\alpha_{2})\frac{s}{(e^{k-sf_{1}}-1)f_{2}+s} & \text{if } 0 \le s < \frac{k}{f_{1}} \\ 1-\alpha_{2} & \text{if } \frac{k}{f_{1}} \le s \le k \end{cases}$$
(2)

within the context of prevalence estimation. The term $f_1 > 0$ defines a threshold for s, namely $\frac{k}{f_1}$, at and above which there is no concentration (or dilution) effect and the term $f_2 > 0$ scales the concentration effect if it is present. Plots of the probability $\pi_1^*(s, k)$ given in formula (2) against s for s = 0, ..., k, k = 8, $\alpha_1 = 0.01$, $\alpha_2 = 0.1$, $f_1 = \frac{4}{3}$, and thus a threshold of s = 6, and selected values of f_2 are presented in Fig. 1(b).

Consider now evaluating the expected number of tests, the expected number of false negatives and the expected number of false positives for an 8 × 12 rectangular array. Suppose that the concentration effect is modeled with $\pi_1^*(s, k)$ of the form (2), that $\alpha_1 = 0.01$ and $\alpha_2 = 0.1$ and that the sensitivity and specificity of individual tests are taken to be $\pi_2^*(+) = 0.9$ and $\pi_2^N(-) = 0.99$. Suppose further, and in contrast to the form for $\pi_1^*(s, k)$ given in Example 1 for which there is no built-in threshold for *s*, that a threshold is introduced here. Specifically, following Hung and Swallow (1999), suppose that $f_1 = 5$ so that the threshold is low and is given by s = 2 for the 8 rows and s = 3 for the 12 columns. Plots of $E_{T,S}(8, 12)$, $E_{FN,S}(8, 12)$ and $E_{FP,S}(8, 12)$ against *p* for S = D, *A*, *O* and *A*2, for $f_2 = 1$ and for a suitable range of prevalences 0 are presented in Fig. 4.

The setting adopted in this example is similar to that of Example 1 but it is clear from a comparison of Figs. 2 and 4 that certain features of the plots for $E_{T,S}(8, 12)$, $E_{FN,S}(8, 12)$ and $E_{FP,S}(8, 12)$ are very different. Thus, while the ordering of the Dorfman, AND, OR and A2 schemes summarized in Table 2 necessarily holds, the shapes of the plots are not the same and the values of the prevalence at which the plotted curves for the schemes cross over are very much higher in the present example than those in Example 1.

These two examples highlight the fact that the operating characteristics for the four schemes of interest depend sensitively on the sub-model adopted for the probability $\pi_1^*(s, k)$, a fact which impacts seriously on the selection of a blood screening protocol. Thus in practice it is usual to have estimates, albeit approximate, for the prevalence, for values of α_1 and α_2 and for the sensitivity and specificity of the individual tests associated with a particular setting (Kim et al., 2007). However selecting a suitable form for the probability $\pi_1^*(s, k)$ is more challenging. Such a selection would most probably necessitate conducting a preliminary experiment involving varying numbers of infected cells in groups of a specified size, followed by the fitting of a suite of sub-models such as those discussed in Examples 1 and 2 to the data in order to find the best-fitting model.

4. Group factor screening

Suppose that a large number of factors that could potentially influence a response are to be investigated in a two-stage group factor screening procedure and that only a small proportion of the factors are effective. Suppose further that each factor can be set at one of two levels, "high" and "low", that a main effects only model is to be invoked to fit the data and that the assumptions for the setting delineated in Watson (1961) for single-group and Morris (2006) for single- and multigroup screening hold. Then the probability that a factor is active, that is effective in influencing the response, is taken to be a constant *p*, the effect of any factor on the response, denoted Δ , is $\Delta > 0$ for all factors that are active and $\Delta = 0$ for all factors not active, and the directions of the effects are known. In the first stage of the group screening procedure the factors are allocated at random to one or more groups and the individual factors within each group are all set at "high" or all set at "low", thereby defining grouped factors at two levels. Suppose here that *g* grouped factors each comprising *k* factors are tested in an appropriate experiment at the α level of significance. Then the probability that a grouped factor is declared active, given that *s* factors within the group are active, where $0 \le s \le k$, can be expressed succinctly as $\pi_1^*(s, k; g; \alpha)$. Note that the dependence of this probability on the number of groups *g* emanates from the nature of the test, and specifically from the form of the test statistic.

Consider now settings for group factor screening which mirror the settings for blood screening discussed in Section 2. Specifically, suppose that the factors are organized in a $k_1 \times k_2$ array and that the Dorfman, AND, OR and A2 schemes are of interest. Then, since k_1 and k_2 are fixed, the arguments relating to the first stage of the procedure are unchanged and the expressions for the expected number of tests are the same, other than the generic notation $\pi_1^*(s, k; g; \alpha)$ which



Fig. 4. Plots of (a) the expected number of tests, (b) the expected number of false negatives and (c) the expected number of false positives against *p* for $0 and the Hung and Swallow (1999) sub-model with <math>k_1 = 8$, $k_2 = 12$, $\alpha_1 = 0.01$, $\alpha_2 = 0.1$, $\pi_2^*(+) = 0.9$, $\pi_2^N(-) = 0.99$, $f_1 = 5$ and $f_2 = 1$.

Y.G. Habtesllassie et al. / Journal of Statistical Planning and Inference I (IIII) III-III

replaces $\pi_1^*(s, k)$. In the second stage of screening the groups are disassembled and all factors within groups which are declared active in the first stage are tested together in an appropriate single experiment. Specifically, suppose that the test $H_0: \Delta = 0$ against $H_A: \Delta > 0$ for an individual factor is performed at the β level of significance. Then the probability that an inactive factor is declared active is β and expressions for the expected number of false positives for the group factor screening schemes of interest are the same as those for the corresponding blood screening schemes. However it follows from the form of the test statistic that the power of the test depends on the number of factors carried through from the first stage, which is a random variable. This dependency must be introduced into the derivation of expressions for the expected number of true positives, and thus of false negatives, as shown for the Dorfman scheme by Watson (1961) and Gurnow (1965).

In the present case, following Gurnow (1965), the expected number of true positives for the AND scheme, is given by

$$E_{TP,A}(k_1, k_2) = \frac{\pi_{R|+}^* \pi_{C|+}^* p}{p \pi_{R|+}^* \pi_{C|+}^* + (1-p) \pi_{R|-}^* \pi_{C|-}^*} \sum_{\substack{\text{feasible}\\n,m}} nm \pi_2^*(+; n, m; \beta) P(N = n, M = m)$$

where *N* and *M* are random variables corresponding to the number of rows and the number of columns declared active in the first stage and $\pi_2^*(+; n, m; \beta)$ denotes the power of the test for N = n and M = m. The expected number of false negatives follows immediately. Note that *N* and *M* are inextricably linked and that their joint distribution would seem to be intractable, at least for large arrays (Habtesllassie, 2004). Thus calculations of $E_{TP,A}(k_1, k_2)$ must necessarily be done by simulation. Similar derivations and considerations hold for the OR and the A2 group factor screening schemes and indeed more broadly for any array-based scheme. It should be emphasized that the results derived here for factor screening are based on the fact that the errors in testing relating to sensitivity in the second stage depend on the numbers of groups testing positive in the first stage and are, as a consequence, in contrast to those for blood screening for which the errors in testing for the two stages are assumed to be independent.

5. Conclusions

In this paper explicit formulae for the expected number of tests, the expected number of false negatives and the expected number of false positives for the Dorfman, the AND, the OR and the A2 group screening schemes for blood samples in the presence of a concentration effect are derived. The arguments used in the derivations are based on those presented by Watson (1961) for group factor screening and Langfeldt et al. (1997) for the blood screening setting. Results for the case of constant sensitivity and specificity in the test procedures and expressions for other operating characteristics, such as pooling positive and negative predictive values, in the array-based schemes of interest follow immediately. In addition it is clear from the illustrative examples that the optimal scheme for a particular operating characteristic depends sensitively on that characteristic, on the prevalence and on the model adopted for the probability $\pi_1^*(s, k)$. In the context of group factor screening, the result for the expected number of false negatives for the Dorfman scheme given in Watson (1961) and Gurnow (1965) is extended to the AND scheme and those for the OR and the A2 schemes can then be derived with minor modification.

There is scope for further research. Thus, in the context of case identification, it would be interesting to find $k_1 \times k_2$ arrays for the four schemes, Dorfman, AND, OR and A2, which comprise the same number of cells and which are in some sense optimal over all feasible integer values of k_1 and k_2 . For example, arrays for which the expected number of tests is a minimum or for which a nominated cost structure is effective could be identified. Some preliminary results in this regard are available in Habtesllassie (2004). Related studies include that for the A2 algorithm with square arrays by Hudgens and Kim (2011) and that for arrays with blockers by Langfeldt et al. (1997). It may also be worthwhile to consider relaxing the assumption that the probability of infection is constant which is made in this study. Specifically, it would be interesting to introduce heterogeneity in the prevalences into the derivations of Section 2.2, following the notions of informative testing discussed in papers by, for example, Bilder et al. (2010) and McMahan et al. (2012). In the context of prevalence estimation, it would seem, from a brief review of the literature, that array schemes such as the AND, the OR and the A2 have not been used in the attendant group screening procedures. It would therefore be worthwhile to investigate whether or not the use of such schemes would impact effectively on the estimation of prevalence, in particular by reducing the mean squared error of the prevalence estimate. Finally, it could well be informative to explore a little further the arguably tenuous link between group screening of blood samples and group factor screening noted in the present study.

Acknowledgments

The authors would like to thank the Associate Editor and two referees for their insightful and helpful comments which did much to improve the paper. The authors would also like to thank the University of Cape Town, the University of KwaZulu-Natal and the National Research Foundation (NRF) of South Africa, grant (UID) 85456, for financial support. The work reported here is based on the Masters thesis of Yonus Habtesllassie. Any opinion, finding and conclusion or recommendation expressed in this material is that of the authors and the NRF does not accept liability in this regard.

14

ARTICLE IN PRESS

Y.G. Habtesllassie et al. / Journal of Statistical Planning and Inference [(IIIII) III-III

Appendix. Derivations for the A2 algorithm

Relating to Section 2.2.4

The expected number of tests associated with the event $R_i^* \cap_{j=1}^{k_2} C_j^N$ follows by first observing that the probability $P\left(R_i^* \cap_{j=1}^{k_2} C_j^N\right)$ can be expressed as the sum of the two terms

$$P\left(R_{i}^{\star}\bigcap_{j=1}^{k_{2}}C_{j}^{\mathsf{N}}|\bigcap_{j=1}^{k_{2}}A_{ij}^{-}\right)P\left(\bigcap_{j=1}^{k_{2}}A_{ij}^{-}\right)$$
(A.1)

and

$$P\left(R_{i}^{\star}\bigcap_{j=1}^{k_{2}}C_{j}^{N}|\text{ at least one }A_{ij}^{+}\right)\times P(\text{ at least one }A_{ij}^{+}).$$
(A.2)

The first term, term (A.1), is given by the product

$$P\left(R_{i}^{\star}|\bigcap_{j=1}^{k_{2}}A_{ij}^{-}\right)P\left(\bigcap_{j=1}^{k_{2}}C_{j}^{N}|\bigcap_{j=1}^{k_{2}}A_{ij}^{-}\right)P\left(\bigcap_{j=1}^{k_{2}}A_{ij}^{-}\right)=\pi_{1}^{\star}(0,k_{2})(\pi_{C|-}^{N}q)^{k_{2}}$$

where $P\left(R_i^{\star} \mid \bigcap_{j=1}^{k_2} A_{ij}^{-}\right) = \pi_1^{\star}(0, k_2)$ and q = 1 - p. The second term, term (A.2), can be written as

$$P(R_i^{\star}|\text{at least one } A_{ij}^+) \times P\left(\bigcap_{j=1}^{k_2} C_j^{N}|\text{at least one } A_{ij}^+\right) \times P(\text{ at least one } A_{ij}^+).$$

Observe now that $P(R_i^*|$ at least one $A_{ij}^+)$ can be interpreted as the average sensitivity over rows of the stage one test procedure and can be expressed as

$$\bar{S}_{(e,R)} = P(R_i^* | \text{at least one } A_{ij}^+) = \frac{\sum_{s=1}^{k_2} \pi_1^*(s, k_2) {\binom{k_2}{s}} p^s (1-p)^{k_2-s}}{(1-q^{k_2})}$$

Also a little reflection shows that

$$P\left(\bigcap_{j=1}^{k_2} C_j^{N} | \text{at least one } A_{ij}^+\right) \times P(\text{ at least one } A_{ij}^+) = P\left(\bigcap_{j=1}^{k_2} C_j^{N} \cap \text{ at least one } A_{ij}^+\right)$$
$$= \left[P\left(\bigcap_{j=1}^{k_2} C_j^{N}\right) - P\left(\bigcap_{j=1}^{k_2} C_j^{N} | \bigcap_{j=1}^{k_2} A_{ij}^-\right) P\left(\bigcap_{j=1}^{k_2} A_{ij}^-\right)\right] = \left[(\pi_C^{N})^{k_2} - (\pi_{C|-}^{N}q)^{k_2}\right].$$

Thus term (A.2) is given succinctly by $\bar{S}_{(e,R)} \left[(\pi_C^N)^{k_2} - (\pi_{C|-}^N q)^{k_2} \right]$. Overall therefore

$$P\left(R_{i}^{\star}\bigcap_{j=1}^{k_{2}}C_{j}^{N}\right) = \left[\pi_{1}^{\star}(0,k_{2}) - \bar{S}_{(e,R)}\right](\pi_{C|-}^{N}q)^{k_{2}} + \bar{S}_{(e,R)}(\pi_{C}^{N})^{k_{2}}$$

It now follows that the probability associated with the event $\bigcap_{i=1}^{k_1} R_i^N \cap C_i^{\star}$ is given by

$$P\left(\bigcap_{i=1}^{k_1} R_i^N \cap C_j^*\right) = \left[\pi_1^*(0, k_1) - \bar{S}_{(e,C)}\right] (\pi_{R|-}^N q)^{k_1} + \bar{S}_{(e,C)} (\pi_R^N)^{k_1}$$

where the term $\bar{S}_{(e,C)}$ is equal to $\bar{S}_{(e,R)}$ but with k_2 replaced by k_1 . The expected number of tests for the $A2(k_1, k_2)$ scheme is thus given by

$$\begin{split} E_{T,A2}(k_1,k_2) &= k_1 + k_2 + k_1 k_2 \left[P(R_i^{\star} \cap C_j^{\star}) + P\left(R_i^{\star} \bigcap_{j=1}^{k_2} C_j^{\mathsf{N}}\right) + P\left(\bigcap_{i=1}^{k_1} R_i^{\mathsf{N}} \cap C_j^{\star}\right) \right] \\ &= k_1 + k_2 + k_1 k_2 \left\{ \pi_{R|+}^{\star} \pi_{C|+}^{\star} p + \pi_{R|-}^{\star} \pi_{C|-}^{\star} (1-p) + \left[\pi_1^{\star}(0,k_2) - \bar{S}_{(e,R)}\right] (\pi_{C|-}^{\mathsf{N}} q)^{k_2} + \bar{S}_{(e,R)} (\pi_C^{\mathsf{N}})^{k_2} + \left[\pi_1^{\star}(0,k_1) - \bar{S}_{(e,C)}\right] (\pi_{R|-}^{\mathsf{N}} q)^{k_1} + \bar{S}_{(e,C)} (\pi_R^{\mathsf{N}})^{k_1} \right\}. \end{split}$$

Y.G. Habtesllassie et al. / Journal of Statistical Planning and Inference & (****) ***-***

The expected number of true positives associated with the event $R_i^* \cap_{j=1}^{k_2} C_j^N$ is evaluated by observing that the conditional probability $P(R_i^* \cap_{j=1}^{k_2} C_j^N \cap A_{ij}^* | A_{ij}^+)$ can be expressed as

$$P(R_i^{\star}|A_{ij}^+)P\left(\bigcap_{\substack{j'=1\\j'\neq j}}^{k_2} C_{j'}^N\right)P(C_j^N|A_{ij}^+)P(A_{ij}^{\star}|A_{ij}^+) = \pi_{R|+}^{\star} (\pi_C^N)^{k_2-1} \pi_{C|+}^N \pi_2^{\star}(+).$$

It then follows immediately that

$$P\left(\bigcap_{i=1}^{k_1} R_i^{\mathsf{N}} \cap C_j^{\star} \cap A_{ij}^{\star} | A_{ij}^+\right) = \pi_{C|+}^{\star} (\pi_R^{\mathsf{N}})^{k_1-1} \pi_{R|+}^{\mathsf{N}} \pi_2^{\star} (+).$$

Overall therefore the expected number of false negatives for the $A2(k_1, k_2)$ scheme is given by

$$\begin{split} E_{FN,A2}(k_1,k_2) &= k_1 k_2 p \left\{ 1 - \left[P(R_i^{\star} \cap C_j^{\star} \cap A_{ij}^{\star} | A_{ij}^+) + P\left(R_i^{\star} \bigcap_{j=1}^{k_2} C_j^{\mathsf{N}} \cap A_{ij}^{\star} | A_{ij}^+\right) + P\left(\bigcap_{i=1}^{k_1} R_i^{\mathsf{N}} \cap C_j^{\star} \cap A_{ij}^{\star} | A_{ij}^+\right) \right] \right\} \\ &= k_1 k_2 p \left\{ 1 - \left[\pi_{R|+}^{\star} \pi_{C|+}^{\star} + \pi_{R|+}^{\star} (\pi_C^{\mathsf{N}})^{k_2 - 1} \pi_{C|+}^{\mathsf{N}} + \pi_{C|+}^{\star} (\pi_R^{\mathsf{N}})^{k_1 - 1} \pi_{R|+}^{\mathsf{N}}] \pi_2^{\star}(+) \right\}. \end{split}$$

Finally, the expected number of false positives associated with the event $R_i^* \cap_{j=1}^{k_2} C_j^N$ can be found by considering the probability $P(R_i^* \cap_{j=1}^{k_2} C_j^N \cap A_{ij}^* | A_{ij}^-) = P(R_i^* \cap_{j=1}^{k_2} C_j^N | A_{ij}^-) P(A_{ij}^* | A_{ij}^-)$. Specifically, the probability $P(R_i^* \cap_{j=1}^{k_2} C_j^N | A_{ij}^-)$ can be found by considering the probabilities associated with the events embedded in the two terms (A.1) and (A.2) conditional on A_{ij}^- . Thus

$$P\left(R_{i}^{\star}\bigcap_{j=1}^{k_{2}}C_{j}^{N}\bigcap_{\substack{j'=1\\j'\neq j}}^{k_{2}}A_{ij'}^{-}|A_{ij}^{-}\right) = P\left(R_{i}^{\star}|\bigcap_{j=1}^{k_{2}}A_{ij}^{-}\right)P\left(\bigcap_{\substack{j'=1\\j'\neq j}}^{k_{2}}(C_{j'}^{N}\cap A_{ij'}^{-})\right)P(C_{j}^{N}|A_{ij}^{-})$$
$$= \pi_{1}^{\star}(0,k_{2})(\pi_{C|-}^{N}q)^{k_{2}-1}\pi_{C|-}^{N}.$$

Also, following the derivation of term (A.2), consider

$$P\left(R_i^{\star}\bigcap_{j=1}^{k_2}C_j^{\mathsf{N}}\cap (\text{at least one }A_{ij'}^+ \text{ with } j'\neq j|A_{ij}^-)\right).$$

Observe first that

$$P(R_i^{\star}|\text{at least one } A_{ij'}^+ \text{ with } j' \neq j \cap A_{ij}^-) = \frac{\sum\limits_{s=1}^{k_2-1} \pi_1^{\star}(s, k_2 - 1) \binom{k_2-1}{s} p^s (1-p)^{k_2-1-s}}{(1-q^{k_2-1})}$$

and, for conciseness, denote this probability by $\bar{S}_{(e,R|-)}$. Then

$$P\left(R_{i}^{*}\bigcap_{j=1}^{k_{2}}C_{j}^{N}\cap(\text{at least one }A_{ij'}^{+}\text{ with }j'\neq j)|A_{ij}^{-}\right)=\bar{S}_{(e,R|-)}\left[(\pi_{C}^{N})^{k_{2}-1}-(\pi_{C|-}^{N}q)^{k_{2}-1}\right]\pi_{C|-}^{N}.$$

Thus

and, similarly,

$$P\left(\bigcap_{i=1}^{k_1} R_i^{\mathsf{N}} \cap C_j^{\star} | A_{ij}^{-}\right) = (\pi_1^{\star}(0, k_1) - \bar{S}_{(e, C|-)})(\pi_{R|-}^{\mathsf{N}} q)^{k_1 - 1} \pi_{R|-}^{\mathsf{N}} + \bar{S}_{(e, C|-)}(\pi_R^{\mathsf{N}})^{k_1 - 1} \pi_{R|-}^{\mathsf{N}}$$

Y.G. Habtesllassie et al. / Journal of Statistical Planning and Inference (()))

where $\bar{S}_{(e,C|-)}$ is given by the expression for $\bar{S}_{(e,R|-)}$ but with *R* replaced by *C* and k_2 by k_1 . Overall therefore the expected number of false positives is given by

$$\begin{split} E_{FP,A2}(k_1,k_2) &= k_1 k_2 \left[P(R_i^{\star} \cap C_j^{\star} | A_{ij}^{-}) + P\left(R_i^{\star} \bigcap_{j=1}^{k_2} C_j^{N} | A_{ij}^{-}\right) + P\left(\bigcap_{i=1}^{k_1} R_i^{N} \cap C_j^{\star} | A_{ij}^{-}\right) \right] \times P(A_{ij}^{\star} | A_{ij}^{-}) P(A_{ij}^{-}) \\ &= k_1 k_2 \left(1-p\right) \left[P(R_i^{\star} \cap C_j^{\star} | A_{ij}^{-}) + P\left(R_i^{\star} \bigcap_{j=1}^{k_2} C_j^{N} | A_{ij}^{-}\right) + P\left(\bigcap_{i=1}^{k_1} R_i^{N} \cap C_j^{\star} | A_{ij}^{-}\right) \right] \pi_2^{\star}(-) \\ &= k_1 k_2 (1-p) \left\{ \pi_{R_i}^{\star} - \pi_{C_i}^{\star} + (\pi_1^{\star}(0,k_2) - \bar{S}_{(e,R_i-)})(\pi_{C_i}^{N} - q)^{k_2 - 1} \pi_{C_i}^{N} + \bar{S}_{(e,R_i-)}(\pi_C^{N})^{k_2 - 1} \pi_{C_i-}^{N} + (\pi_1^{\star}(0,k_1) - \bar{S}_{(e,C_i-)})(\pi_{R_i}^{N} - q)^{k_1 - 1} \pi_{R_i-}^{N} + \bar{S}_{(e,C_i-)}(\pi_R^{N})^{k_1 - 1} \pi_{R_i-}^{N} \right\} \pi_2^{\star}(-). \end{split}$$

Relating to Section 2.2.5

Observe that for sensitivities and specificities which do not depend on the number of infected cells in a group

$$P\left(R_{i}^{\star}\bigcap_{j=1}^{k_{2}}C_{j}^{N}\right) = (1-S_{p}-S_{e})(1-p)^{k_{2}}\left[1-t(k_{1}-1)\right]^{k_{2}} + S_{e}\left[1-t(k_{1}-1)\right]^{k_{2}}$$

that

$$P\left(R_{i}^{\star}\bigcap_{j=1}^{k_{2}}C_{j}^{N}\cap A_{ij}^{\star}|A_{ij}^{+}\right)=S_{e}^{2}(1-S_{e})\left[1-t(k_{1})\right]^{k_{2}-1}$$

and that

$$P\left(R_{i}^{\star}\bigcap_{j=1}^{k_{2}}C_{j}^{N}\cap A_{ij}^{\star}|A_{ij}^{-}\right) = (1-S_{p})\left\{(1-S_{p}-S_{e})\left(1-p\right)^{k_{2}-1}\left[1-t(k_{1}-1)\right]^{k_{2}}+S_{e}\left[1-t(k_{1})\right]^{k_{2}-1}\left[1-t(k_{1}-1)\right]\right\}$$

for $i = 1, ..., k_1$. Probabilities associated with the events $\bigcap_{j=1}^{k_1} R_i^N \cap C_j^*$, $j = 1, ..., k_2$, follow immediately by interchanging k_1 and k_2 in the above expressions. Note that the sums embedded in the terms h(n) and h(n|y) given in the paper by Kim et al. (2007) can be evaluated explicitly and thus that the results presented here are in accord with the expressions for $P(R_i^* \cap_{j=1}^n C_j^N)$, the pooling sensitivity and the pooling specificity developed for the A2(n, n) scheme in that paper. The expected number of tests, the expected number of false negatives and the expected number of false positives for the $A2(k_1, k_2)$ scheme now follow immediately from the general formulae given in Section 2.2.4 and are given by

$$\begin{split} E_{T,A2}(k_1,k_2) &= k_1 + k_2 + k_1 k_2 \left\{ S_e^2 \ p + t(k_1 - 1)t(k_2 - 1) \ (1 - p) \right. \\ &+ (1 - S_e - S_p) \left\{ \left[(1 - t(k_1 - 1)) \ q \right]^{k_2} + \left[(1 - t(k_2 - 1)) \ q \right]^{k_1} \right\} \\ &+ S_e \left\{ \left[(1 - t(k_1 - 1)) \right]^{k_2} + S_e \left[(1 - t(k_2 - 1)) \right]^{k_1} \right\} \right\} , \\ E_{FN,A2}(k_1,k_2) &= k_1 k_2 \ p \ \left\{ 1 - S_e^3 - S_e^2 (1 - S_e) \left[(1 - t(k_1))^{k_2 - 1} + (1 - t(k_2))^{k_1 - 1} \right] \right\} \end{split}$$

and

$$\begin{split} E_{FP,A2}(k_1,k_2) &= k_1 k_2 (1-p) (1-S_p) \left\{ t (k_1-1) t (k_2-1) \right. \\ &+ (1-S_e-S_p) \left\{ [1-t (k_1-1)]^{k_2} q^{k_2-1} + [1-t (k_2-1)]^{k_1} q^{k_1-1} \right\} \\ &+ S_e \left\{ [1-t (k_1)]^{k_2} + [1-t (k_2)]^{k_1} \right\} \end{split}$$

respectively.

References

Bilder, C.R., Tebbs, J.M., Chen, P., 2010. Informative retesting. J. Amer. Statist. Assoc. 105, 942–955.

Burns, K.C., Mauro, C.A., 1987. Group testing with test error as a function of concentration. Comm. Statist. Theory Methods 16, 2957–2979.

Dean, A.M., Lewis, S.M., 2006. Screening Methods for Experimentation in Industry, Drug Discovery and Genetics. Springer, New York.

Dorfman, R., 1943. The detection of defective members of large populations. Ann. Math. Stat. 14, 436-440.

Gurnow, R., 1965. A note on G.S. Watson's paper 'A study of the group screening method'. Technometrics 7, 444-446.

Habtesllassie, Y.C., 2004. Group Screening with Imperfect Testing (Master's thesis), University of Natal, Pietermaritzburg, South Africa. Hedt, B.L., Pagano, M., 2008a. A Matrix Pooling Algorithm for Disease Detection. Tech. Rep. 57. Harvard University Biostatistics, Working Paper Series. Hedt, B.L., Pagano, M., 2008b. Matrix Pooling: An Accurate and Cost Effective Testing Algorithm for Detection of Acute HIV Infection. Tech. Rep. 58. Harvard University Biostatistics Working Paper Series.

Hudgens, M.G., Kim, H.-Y., 2011. Optimal onfiguration of a square array group testing algorithm. Comm. Statist. Theory Methods 40, 436-448.

Hughes-Oliver, I.M., 2006. Pooling experiments for blood screening and drug discovery. In: Dean, A.M., Lewis, S.M. (Eds.), Screening Methods for Experimentation in Industry, Drug Discovery and Genetics. Springer, New York, pp. 48-68.

Please cite this article in press as: Habtesllassie, Y.G., et al., Array-based schemes for group screening with test errors which incorporate a concentration effect. J. Statist. Plann. Inference (2015), http://dx.doi.org/10.1016/j.jspi.2015.05.009

16

Y.G. Habtesllassie et al. / Journal of Statistical Planning and Inference I (IIII) III-III

Hung, M., Swallow, W.H., 1999. Robustness of group testing in the estimation of proportions. Biometrics 55, 231–237.

Hwang, F.K., 1976. Group testing with a dilution effect. Biometrika 63, 671–680.

Kim, H.-Y., Hudgens, M.G., 2009. Three-dimensional array-based group testing algorithms. Biometrics 65, 903–910.

Kim, H.-Y., Hudgens, M.G., Dreyfuss, J.M., Westreich, D.J., Pilcher, C.D., 2007. Comparison of group testing algorithms for case identification in the presence of test error. Biometrics 63, 1152–1163.

Langfeldt, S.A., Hughes-Oliver, J.M., Ghosh, S.K., Young, S.S., 1997, Optimal Group Testing in the Presence of Blockers. Technical Report, Institute of Statistics Mimeograph Series No. 2297, North Carolina University, Raleigh NC, USA.

Liu, A., Liu, Č., Zhang, Z., Albert, P.S., 2012. Optimality of group testing in the presence of misclassification. Biometrika 99, 245–251.

McMahan, C.S., Tebbs, J.M., Bilder, C.R., 2012. Informative Dorfman screening. Biometrics 68, 287-296.

Morris, M.D., 1987. Two-stage factor screening procedures using multiple grouping assignments. Comm. Statist. Theory Methods 16, 3051–3067. Morris, M.D., 2006. An overview of group factor screening. In: Dean, A.M., Lewis, S.M. (Eds.), Screening Methods for Experimentation in Industry, Drug

Discovery and Genetics. Springer, New York, pp. 191–206.

Phatarfod, R.M., Sudbury, A., 1994. The use of a square array scheme in blood testing. Stat. Med. 13, 2337–2343.

Watson, G.S., 1961. A study of the group screening method. Technometrics 3, 371–388.

Wein, L.M., Zenios, S.A., 1996. Pooled testing for hiv screening: Capturing the dilution effect. Oper. Res. 44, 543–569.

Zhang, Z., Liu, C., Kim, S., Liu, A., 2014. Prevalence estimation subject to misclassification: the mis-substitution bias and some remedies. Stat. Med. 33, 4482-4500.