

Abstract

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24

Background: *The sample size required to power a study to a nominal level in a paired comparative diagnostic accuracy study, i.e. studies in which the diagnostic accuracy of two testing procedures is compared relative to a gold standard, depends on the conditional dependence between the two tests - the lower the dependence the greater the sample size required. A priori, we usually do not know the dependence between the two tests and thus cannot determine the exact sample size required. One option is to use the implied sample size for the maximal negative dependence, giving the largest possible sample size. However, this is potentially wasteful of resources and unnecessarily burdensome on study participants as the study is likely to be overpowered. A more accurate estimate of the sample size can be determined at a planned interim analysis point where the sample size is re-estimated.*

Methods: *This paper discusses a sample size estimation and re-estimation method based on the maximum likelihood estimates, under an implied multinomial model, of the observed values of conditional dependence between the two tests and, if required, prevalence, at a planned interim. The method is illustrated by comparing the accuracy of two procedures for the detection of pancreatic cancer, one procedure using the standard battery of tests, and the other using the standard battery with the addition of a PET/CT scan all relative to the gold standard of a cell biopsy. Simulation of the proposed method illustrates its robustness under various conditions.*

Results: *The results show that the type I error rate of the overall experiment is stable using our suggested method and that the type II error rate is close to or above nominal. Furthermore, the instances in which the type II error rate is above nominal are in the situations where the lowest sample size is required, meaning a lower impact on the actual number of participants recruited.*

1 **Conclusion:** *We recommend multinomial model maximum likelihood estimation of the*
2 *conditional dependence between paired diagnostic accuracy tests at an interim to reduce the*
3 *number of participants required to power the study to at least the nominal level.*

4 *Trial registration: ISRCTN ISRCTN73852054. Registered 9th of January 2015.*

5 *Retrospectively registered.*

6

7 **Keywords:** Interim analysis, sample-size re-estimation, study design, diagnostic accuracy,
8 sensitivity, specificity.

9

10

11

12

13

14

15

16

17

18

19

20

1 **Background**

2 An assessment of diagnostic accuracy is crucial in the development of medical testing
3 procedures [1]. Comparing the accuracy of these procedures in terms of their sensitivities
4 and specificities [2,3] relative to a gold standard, a type of ‘diagnostic accuracy’ study [4], is
5 essential to ensuring that the most appropriate tests are deployed in the clinical setting [5]. At
6 the outset of a study, a sample size is calculated based on assumptions made about the
7 expected changes in sensitivity and specificity and, in a prospective design, the likely
8 prevalence of the condition to be tested for in the sample. However, the initial assumptions
9 about parameters in the study, especially the conditional dependence between the two tests,
10 may be revealed to be inaccurate, resulting in a potentially over- or under-powered study. A
11 planned interim analysis can allow the study's sample size to be updated based on the data
12 already collected. This involves utilising the information observed at the interim stage to
13 refine the sample size estimate. A resulting increase in sample size allows the time, cost and
14 patient discomfort already invested in the study to yield valid results while a decrease in
15 sample size means that less time and cost will be expended overall and patients will not
16 needlessly undergo unnecessary testing [6].

17 There are well-established methodologies for interim sample size re-estimation in treatment
18 studies for continuous and normally distributed response variables [7–11], some of which
19 provide mechanisms to maintain blinding in the study [8–10]. Methods also exist for the re-
20 estimation with binary response variables [12,13], and mechanisms to maintain blinding have
21 been proposed in this more complex situation where the variance and mean parameters are
22 not separable [14]. Proschan [15] gives an overview of sample size re-estimation procedures
23 based on a nuisance parameter. Specifically, procedures for determining the difference of
24 means between two samples with a common, unknown, variance and difference in

1 proportions between two groups, with an unknown overall proportion, are considered. In the
2 case of normally distributed data, the independence of the sample variance and sample mean
3 ensures that the validity of estimates is unaffected by the interim sample size re-estimation
4 and this is shown to hold asymptotically in the binary case. However, Proschan does not
5 consider the case of paired data which is the focus of the current paper. Furthermore, the
6 implications of sample size re-estimation in the context of comparative diagnostics studies,
7 inherently different from those in treatment (randomised controlled) studies [16], have not
8 been fully explored in the statistical literature.

9 A number of salient differences in interim analysis between studies comparing diagnostic
10 tests and those comparing treatments are highlighted in Gerke *et al* [5] and Gerke *et al* [16].
11 Firstly, in paired diagnostics accuracy studies, full blinding is often not possible. However,
12 as long as the results of the two-tests which are being compared are temporarily blinded from
13 each other, this is not a major threat to a study's validity. In fact, it has the advantage that the
14 patients can benefit from their clinicians knowing the results of both diagnostic tests after
15 testing has taken place. Secondly, in diagnostic accuracy studies, early cessation of the study
16 due to futility is not as easy to establish as in treatment studies. The reason for this is the fact
17 that treatment studies often test a single outcome while diagnostic studies test two outcomes,
18 sensitivity and specificity, and futility must be established for both simultaneously. Thirdly,
19 the sample size required for a hypothesis test in diagnostic studies, powered to a given level,
20 is closely related to the conditional dependence between the two testing procedures which has
21 been shown to present problems in a number of contexts [5,17–22]. More specifically, the
22 lower the conditional dependence between the tests, the greater the sample size will be, with
23 the largest sample size being implied by the maximum negative dependence, given the
24 specified alternative hypotheses. This level of conditional dependence between the tests is
25 one of the primary factors driving the required sample size estimate and it is often difficult to

1 estimate *a priori*. Gerke *et al* [5] assert that for comparative diagnostic studies, as long as an
2 interim sample size re-estimation is planned it bears no threat to the validity of the study.
3 However, Gerke *et al* [5] do not provide justification for this assertion and, furthermore, their
4 assertion does not take the inherent uncertainty of the interim data into account. This study
5 aims to present a method and give practical guidelines for its application, for the initial
6 estimation and interim re-estimation of sample size in a paired diagnostic study which will
7 allow utilisation of information on the conditional dependence between tests at the interim to
8 potentially reduce the required sample size while maintaining the approximate nominal
9 statistical power of the experiment as a whole. While we present a method of estimating the
10 size of the conditional dependence to reduce sample size, it should also be noted that there is
11 a body of literature dealing with the problems caused by conditional dependence in other
12 areas [23–25].

13 The remainder of the article is organised as follows. The methods section outlines sample
14 size estimation methods for paired diagnostic test studies, introduces a motivating example
15 application, and then goes on to propose a new method for re-estimation based on a
16 multinomial likelihood. The results section first provides extensive simulations of the method
17 under various real world conditions and then moves back to apply the sample size re-
18 estimation method proposed in this paper to the motivating example. The article then
19 continues with a brief discussion of the place of this study in the literature and the optimal
20 interim sample size to choose. Finally, the conclusion, summarises and restates the major
21 outcomes of this study.

22 **Methods**

23 A representation of data from a paired comparative diagnostic accuracy study is given in
24 Table 1. The subjects are initially divided according to whether they are discovered, via the

1 gold standard test, to be diseased or non-diseased. They are then further subdivided as to
 2 whether they test positive or negative on tests A and B. For example, the cell n_A represents
 3 subjects that were found to have the disease via the gold standard test and also tested positive
 4 on both test A and B, while cell n_F denotes subjects who tested negative on the gold standard
 5 and test B but positive on test A.

6 **Table 1- Paired study design**

		Diseased		Non-diseased			
		Test B		Test B			
		+ive	-ive	+ive		-ive	
Test A	+ive	n_A	n_B	Test A	+ive	n_E	n_F
	-ive	n_C	n_D		-ive	n_G	n_H

7
 8 A possible initial sample size calculation, using a normal approximation of the logarithm of
 9 the ratio of sensitivities and specificities, and assuming a comparison between a new test, test
 10 A, and an existing test, test B, follows from Alonzo *et al* [19] and a full derivation can be
 11 found therein. The experiment, as a whole tests jointly both sensitivity and specificity
 12 improvement to pre-specified levels, the sample size is calculated for each and the largest
 13 sample size is chosen to power the study. Note that this paper concentrates on the situation in
 14 which superiority is tested in for both sensitivity and specificity. However, the method
 15 elaborated below should be extendable to situations where we are interested in testing non-
 16 inferiority in either or both of sensitivity and specificity. For details on the construction of
 17 the confidence intervals and hypothesis tests in these situations see Alonzo *et al* [19]. In the
 18 case of the estimation of a sample size for superiority, the initial sample size calculation for
 19 sensitivity is given by:

$$n_{p1} = \left(\frac{Z^{(1-\beta)} + Z^{(1-\alpha/2)}}{\log \gamma_1} \right)^2 \left(\frac{(\gamma_1 + 1)TPR_B - 2TPPR}{\gamma_1 TPR_B^2} \right) / \pi \quad (1)$$

1 where, α is the type I error rate of the study and β is the power of the study. The main
2 quantity of interest, γ_1 , is the ratio of true positive rates = TPR_A/TPR_B , TPR_B is the true
3 positive rate (sensitivity) on test B, i.e. $TPR_B = (n_A + n_C) / (n_A + n_B + n_C + n_D)$, TPR_A
4 is the true positive rate (sensitivity) on test A, i.e. $TPR_A = (n_A + n_B) / (n_A + n_B + n_C +$
5 $n_D)$, $TPPR$ is the proportion of diseased patients who test positive on both tests, i.e.
6 $TPPR = n_A / (n_A + n_B + n_C + n_D)$ and π is the prevalence of disease. The null hypothesis
7 is that $\gamma_1 = 1$, the alternative hypothesis is that $\gamma_1 \neq 1$.

8 For testing superiority of specificity we are interested in the true negative rates so the formula
9 is instead:

10

$$n_{n1} = \left(\frac{Z^{(1-\beta)} + Z^{(1-\alpha/2)}}{\log \gamma_2} \right)^2 \left(\frac{(\gamma_2 + 1)TNR_B - 2TNNR}{\gamma_2 TNR_B^2} \right) / (1 - \pi) \quad (2)$$

11 where, γ_2 , the main quantity of interest is the ratio of true negative rates = TNR_A/TNR_B ,
12 TNR_A is the true negative rate (specificity) on test A = $(n_G + n_H) / (n_E + n_F + n_G + n_H)$,
13 TNR_B is the true negative rate (specificity) on test B = $(n_F + n_H) / (n_E + n_F + n_G + n_H)$,
14 and $TNNR$ is the proportion of non-diseased patients who test negative on both tests =
15 $n_H / (n_E + n_F + n_G + n_H)$.

16 It is interesting to note that, following the notation of Vacek [23] and considering the
17 population 2x2 table (in Table 1), the conditional dependence of the two tests can be denoted
18 by e_b and e_a , the conditional covariance when the gold standard disease status is positive or
19 negative, respectively [23]. Therefore, the probability of both tests being positive can be

1 expressed as $TPPR = TPR_A \cdot TPR_B + e_b$ and the probability of both tests being negative
2 $TNNR = (1 - TNR_A) \cdot (1 - TNR_B) + e_a$. When e_a and $e_b = 0$ the tests are conditionally
3 independent, when e_a and/or $e_b \neq 0$ the response on one test changes the probability of that
4 response on the other test. For example, when $e_b > 0$ an individual who responds positively
5 on test A is more likely to respond positively on test B.

6

7 For initial estimates of $TPPR$ and $TNNR$, from Alonzo *et al* [19] we can use the fact that
8 $TPPR \geq (1 + \gamma_1)TPR_B - 1$ and $TNNR \geq (1 + \gamma_2)TNR_B - 1$ to estimate the lower bounds
9 of the possible values of $TPPR$ and $TNNR$, under the specified hypotheses. The required
10 sample size is largest when $TPPR = (1 + \gamma_1)TPR_B - 1$ and $TNNR = (1 + \gamma_2)TNR_B - 1$,
11 thus, these estimates represent the "worst case scenarios" of maximal negative conditional
12 dependence between the tests, conditional on the fixed values of TPR_A and TPR_B . The
13 sample size implied by using these levels of $TPPR$ and $TNNR$ would very likely overpower
14 the study, i.e. more participants will be recruited than is strictly necessary to achieve the
15 power specified by β . The required sample size is smallest when the conditional dependence
16 between tests A and B are maximal, conditional on the fixed values of TPR_A and TPR_B , i.e.
17 when $TPPR = TPR_B$ and $TNNR = TNR_B$. The implied sample size in this case would likely
18 underpower the study, i.e. too few participants recruited to reach the power specified by β .
19 The sample size in this "best case scenario" can be substantially lower than that in the worst
20 case scenario. Conservatively, it might be thought a good idea to always use the "worst case
21 scenario" implied sample size estimate which will always power the study sufficiently.
22 However, in cases where the recruitment and testing of participants comes at a premium, both
23 financially and in terms of discomfort to the patients, it might be preferable to apply a more
24 nuanced strategy. Furthermore, the sample size implied by the "worst case scenario" implies

1 the highly unlikely condition of a maximal negative conditional dependence between two
2 tests, which are performed on the same patients to detect the same disease. The implied
3 sample size based on this condition is not recommended [26]. One possibility, to enable a
4 more accurate evaluation of the conditional dependence between the two tests, and thus the
5 required sample size, is to perform a planned interim sample size re-estimation using this
6 information to refine the sample size estimate.

7

8 At a planned interim, where a proportion of the overall sample size has been collected, we
9 would have some information about the true values of $TPPR$, $TNNR$, π , TPR_B and TNR_B ,
10 however, these values would only come from a limited sample size. The crucial parameters
11 to use in re-estimation are those related to the conditional dependence between the tests, i.e.,
12 $TPPR$ and $TNNR$, as these values are difficult to estimate and, for these parameters, it is
13 unlikely that research exists which can provide an approximate value. Conversely, the values
14 of, TPR_B and TNR_B , the sensitivities and specificities of an established test, may have known
15 values in the literature and these should preferably be used over those from the relatively
16 small interim sample. For the value of π , the prevalence, a judgement must be made as to
17 whether the researcher feels that any pre-existing estimate of prevalence would be a more
18 accurate reflection of the true prevalence in the specific study population than any interim
19 estimate. In the example given below, we use values for $TPPR$, $TNNR$ and π at the interim
20 in the sample size calculation.

21 Naively, it might appear that interim sample size re-estimation would entail a straightforward
22 replication of equations (1) and (2) with π , and in the case of (1), $TPPR$ or in the case of (2),
23 $TNNR$, replaced with the estimates at the interim point. However, this approach does not
24 effectively take into account the inherent uncertainty in the interim parameter estimates of

1 $TPPR$, $TNNR$ and π , nor the fact that only a specific range of values for $TPPR$ and $TNNR$
2 are actually possible under the alternative hypothesis. An approach which does take these
3 factors into account is re-estimation of the sample size based on maximum likelihood
4 estimation, at the interim, of the parameters in question under a multinomial model. This
5 model is constrained by the hypothesised values of TPR_A, TPR_B , TNR_A , and TNR_B , i.e. the
6 marginals in Table 1.

7

8 **Application**

9 The numerical example we use involves an interim sample size recalculation of a study
10 comparing the incremental benefits to sensitivity and specificity of augmenting current
11 methods for diagnosing pancreatic cancer with Positron Emission Tomography (PET) and
12 computed tomography (CT) technologies. The alternative hypotheses were that sensitivity
13 would rise from 81% to 90%, and specificity would rise from 66% to 80%, additionally, the
14 expected prevalence of pancreatic cancer from the literature was 47%.

15

16 To calculate the sample size for sensitivity equation 1 was used, taking $\alpha = 0.05$, $\beta = 0.2$,
17 $\hat{\gamma}_1 = \frac{0.9}{0.81}$, $\widehat{TPR}_B = 0.81$, $\widehat{TPPR} = 0.71$, and $\hat{\pi} = 0.47$ gives a sample size of **598**. To
18 calculate the sample size for specificity equation 2 was used taking $\alpha = 0.05$, $\beta = 0.2$, $\hat{\gamma}_2 =$
19 $\frac{0.8}{0.66}$, $\widehat{TNR}_B = 0.66$, $\widehat{TNNR} = 0.46$, and $\hat{\pi} = 0.47$ gives a sample size of **409**. The minimum
20 sample sizes for sensitivity and specificity, given $\widehat{TPPR} = 0.81$ and $\widehat{TNNR} = 0.66$, are **186**
21 and **106**, respectively. Given the disparity between the minimum and maximum sample size
22 estimates it was decided to re-assess the sample size at a planned interim.

23

Table 6 - Interim PET diagnostic study results

Diseased patients				Non-diseased patients			
		Pre-PET				Pre-PET	
		+ive	-ive			+ive	-ive
Post-PET	+ive	66	3	Post-PET	+ive	21	4
	-ive	3	10		-ive	11	69

1
2 Table 6 gives the results after data from 187 participants had been collected. The observed
3 values at the interim are: $\widehat{TPPR} = 0.80$, $\widehat{TNNR} = 0.66$ and $\hat{\pi} = 0.44$. Taking a naive
4 approach and plugging these values directly into equations 1 and 2 the implied sample sizes
5 for sensitivity become **242** and for specificity **100**, giving a total sample size for the study of
6 **242** (or **342** and **145**, respectively, had we also used the interim values of TPR_B and TNR_B).
7 However, this method does not take into account the fact that \widehat{TPPR} and \widehat{TNNR} are random
8 variables and we are actually interested in the true value of the probability of $TPPR$ and
9 $TNNR$ under the specified alternative hypothesis. In fact, had the observed value for $TPPR$
10 been equal to 0.86, the sample size given via the naive method would have been **-22**, given
11 the fact that \widehat{TPPR} would have been larger than both TPR_A and TPR_B . Clearly, the naive
12 method, which uses the random value of a single cell, is inappropriate and a method that uses
13 information about the value of $TPPR$ from all of the observed cells and the specified
14 marginals is required.

15

16 **Sample size re-estimation via maximum likelihood estimation of $TPPR$**

17 For illustration purposes, we will discuss the re-estimation of the sample size for sensitivity,
18 the estimation procedure for specificity is analogous. Taking TPR_A as the test with the
19 highest expected diagnostic utility, i.e. the “new” test whose performance we are comparing
20 to the “standard”, the probabilities corresponding to the cells in Table 1, given the situation of
21 the maximally negative conditional dependence between the tests are: $p_1 = TPR_B - (1 -$

1 TPR_A), $p_2 = 1 - TPR_B$, $p_3 = 1 - TPR_A$, $p_4 = 0$. The probabilities of the cells when the
 2 conditional dependence between TPR_A and TPR_B is at its maximally positive are given
 3 by: $p_1 = TPR_B$, $p_2 = TPR_A - TPR_B$, $p_3 = 0$, $p_4 = 1 - TPR_A$. We could alternatively
 4 specify these cell probabilities according to the covariance between the two tests.
 5 Specifically, Vacek [23] gives the maximum value of the covariance as $TPR_B (1 -$
 6 $TPR_A)$ and the minimum value as $-(1 - TPR_A)(1 - TPR_B)$. Thus, the maximum and
 7 minimum values for the cells can be ascertained by finding the product of the marginal
 8 probabilities associated with a cell and adding the minimum or maximum value of
 9 covariance, for cells p_1 and p_4 , or subtracting the values of covariance for cells p_2 and p_3 .
 10 For example, the minimum value for $p_1 = TPR_A \cdot TPR_B - (1 - TPR_A)(1 - TPR_B)$.
 11 Between the minimum and maximum values lies every permissible joint configuration. Let
 12 these possible joint configurations be expressed as vector, \mathbf{p} , with $p_1 = TPR$, where
 13 $\sum_{i=1}^4 p_i = 1$, $p_1 + p_2 = TPR_A$ and $p_1 + p_3 = TPR_B$.

14

15 When the conditional dependence is maximally positive the sample size required is the
 16 smallest, when it is maximally negative the sample size required is at its largest. At the
 17 beginning of the experiment we do not know which of these possible levels of conditional
 18 dependence our data were generated under and thus we use the, usually overly conservative,
 19 largest possible sample size estimate.

20

21 However, at the interim we can use our observed data to infer a likelihood of that data having
 22 been generated under each of the permissible joint configurations of cell probabilities given
 23 the implied range of probabilities under a multinomial model. A simple method of

1 extracting an estimate of TPPR is to maximise the likelihood function of the interim data
 2 given the values of \mathbf{p} implied by the marginal probabilities:

$$\mathcal{L}(\mathbf{p}|x) = \prod_{i=1}^4 \mathbf{p}_i^{x_i} \quad (3)$$

3 where \mathbf{p} is the vector of joint probabilities defined above and x are the observed cell
 4 frequencies. The constraints imposed on the above multinomial likelihood make the
 5 parameter space one dimensional, thus, substituting the constraints in order to express the
 6 likelihood in terms of p_1 , gives:

$$\mathcal{L}(p_1|x) = p_1^{x_1} (TPR_A - p_1)^{x_2} (TPR_B - p_1)^{x_3} (1 - TPR_A - TPR_B + p_1)^{x_4} \quad (4)$$

$$p_1 \in [TPR_B - (1 - TPR_A), TPR_B]$$

7 Code to estimate this in R, via optimisation of the negative log-likelihood, is in the Appendix.
 8 In effect, this method bounds the value for the conditional dependence between the minimum
 9 and maximum values under the specified marginals and then uses information from the
 10 frequency values of the four cells of the table to infer the most probable value of p_1 . We can
 11 use this estimate of \hat{p}_1 as our value of \widehat{TPPR} and use the observed value of the prevalence (if
 12 required) as our measure of $\hat{\pi}$ in equation 1 to re-estimate the sample size at the interim.

13 **Results**

14 **Simulation Studies**

15 In order to verify the integrity of the method for sample size re-estimation described and
 16 applied above a series of simulation studies were carried out. The objectives of these studies
 17 were to assess the implications of re-estimating a sample size based on data already collected
 18 on the type I and II error rates under various permutations of parameters. The type II error

1 rate should be as close to nominal as possible (i.e. 0.8 in the example above), and the type I
2 error rate should be minimally affected by the re-estimation.

3 It should be noted that the statistical power provided by the sample size implied by the
4 Alonzo et al [19] method (when no re-estimation is undertaken) is related to the level of
5 conditional dependence between the tests, Figure 1 illustrates this relationship. In total
6 100,000 replications were generated under the specified true alternative hypothesis (i.e. $\gamma_1 =$
7 $0.9/0.81 = 1.11$), for the example situation above, at various levels of conditional dependence
8 between the two tests. The number of replications 100,000, is more than required, however
9 as the computing time to calculate these was trivial, there was little cost in simulating to this
10 level of accuracy. This number of simulations was used throughout this paper. In all cases in
11 Figure 1 the simulated power was higher than nominal but where the conditional dependence
12 was highest the power was greatly over specified. As the conditional dependence tends
13 towards becoming maximally positive, i.e. as TPPR tends towards its maximal value, the cell
14 n_C tends towards 0. This means that the asymptotic assumptions underlying formulae 1 and
15 2 and those underlying the significance test no longer hold. However, this should not be of
16 too great a concern, with regards to balancing the minimisation of the required sample size
17 estimate with the statistical power of the experiment, as the instances where the power is over
18 specified are when the sample size is lowest. Additional conservatism at positive levels of
19 conditional dependence has a significantly lesser impact on the overall sample size than it
20 would have at the end of the continuum where the conditional dependence is negative.
21 Whatever the case may be, it should be noted that the results of re-estimation will follow a
22 similar pattern.

23 **Figure 1** – Simulated power of sample size specified by the true TPPR in equation 1 when
24 $TPR_A=0.9$, $TPR_B=0.81$ and $\pi=0.45$.

1

2

Figure 1 here

3

In the first set of simulations, which aim to assess the stability of the type II error rate, data

4

are generated under the conditions $TPR_A = 0.9$, $TPR_B = 0.81$, $\pi = 0.45$, while the sample

5

sizes at the interim are varied between 50 and 200 and the values for $TPPR$ are varied

6

between 0.71 and 0.81. The null hypothesis is: $TPR_A/TPR_B = 1$, and our data were simulated

7

under the alternative hypothesis $TPR_A = 0.9$ and $TPR_B = 0.81$, with varying levels of

8

conditional dependence within the implied limits. Figure 2 shows how the power of the

9

experiment overall (i.e. using the data from both before and after sample size re-estimation)

10

varies as a function of the interim sample size and the true value of $TPPR$. As expected the

11

values follow the same pattern as that in Figure 1. The minimum of the nominal power, or

12

very close to it, was achieved at all levels of conditional dependence and at all interim sample

13

sizes.

14

Figure 2 – Simulated power of re-estimation method across various interim sample sizes and

15

levels of true $TPPR$ when $TPR_A = 0.9$, $TPR_B = 0.81$ and $\pi = 0.45$.

16

17

Figure 2 here

18

Table 2 provides information about the mean sample size, bias, coverage and Root Mean

19

Squared Error (RMSE) (from the value specified by equation 1 using the true value of $TPPR$

20

for the simulated data) under the combinations of conditional dependences and interim

21

sample size. The sample sizes implied by Equation 1 for maximal and minimal levels of

22

conditional dependence are **194** and **625**, respectively. The interim sample sizes of 50, 100,

23

150 and 200 were chosen to illustrate the effects of choosing various interim sample sizes

1 that were smaller than the total sample size of 242 calculated by the Alonzo method
 2 described above for our application.

3 An increasing interim sample size does not have that great an impact on the average
 4 estimated sample size. However, it does have a large impact on the RMSE. Thus, choosing a
 5 larger interim sample size at which to re-estimate will ensure a more accurate sample size re-
 6 estimate in individual cases, meaning that the experiment will be more likely to be powered
 7 to the appropriate level while recruiting as few participants as possible. Of course, if the
 8 interim sample size is chosen to be too large then there is a risk of having already recruited
 9 too many participants at the interim. Therefore, some sensible trade-off is required. The bias
 10 and coverage seem to be at acceptable levels although the coverage does dip when the
 11 conditional dependence between the tests is high.

12 **Table 2** – Mean sample size (S.D.), bias, coverage and RMSE of simulated sample sizes with
 13 varying interim sample size estimates and true levels of TPPR when $TPR_A = 0.9$, $TPR_B =$
 14 0.81 and Prevalence = 0.45. (N=interim sample size.)

	<u>Mean sample size</u>				<u>Bias</u>			
	<u>N=50</u>	<u>N=100</u>	<u>N=150</u>	<u>N=200</u>	<u>N=50</u>	<u>N=100</u>	<u>N=150</u>	<u>N=200</u>
TPPR								
0.81	217(77)	202(35)	198(23)	205(17)	-0.00091	-0.00027	0.00018	0.00048
0.80	256(114)	241(71)	238(56)	241(48)	-0.00031	0.00035	0.00062	0.00064
0.79	297(139)	283(92)	281(72)	282(62)	-0.00007	0.00069	0.00072	0.00068
0.78	338(155)	326(105)	325(83)	325(70)	0.00045	0.00056	0.00082	0.00062
0.77	381(166)	371(114)	369(89)	369(75)	0.00043	0.00054	0.00058	0.00050
0.76	423(170)	415(118)	413(92)	413(78)	0.00054	0.00035	0.00054	0.00041
0.75	465(171)	460(118)	457(93)	456(79)	0.00069	0.00056	0.00029	0.00033
0.74	506(166)	503(115)	501(91)	500(78)	0.00029	0.00028	0.00031	0.00031
0.73	546(156)	546(107)	545(86)	543(73)	0.00047	0.00045	0.00022	0.00022
0.72	585(143)	588(95)	88(76)	586(65)	0.00043	0.00027	0.00017	0.00022
0.71	621(124)	629(75)	630(59)	629(50)	0.00024	0.00037	0.00033	0.00019

	<u>Coverage</u>				<u>RMSE</u>			
	<u>N=50</u>	<u>N=100</u>	<u>N=150</u>	<u>N=200</u>	<u>N=50</u>	<u>N=100</u>	<u>N=150</u>	<u>N=200</u>
TPPR								
0.81	0.923	0.925	0.924	0.923	80	36	23	18
0.8	0.936	0.937	0.936	0.936	115	71	62	48
0.79	0.942	0.943	0.944	0.943	140	92	72	62
0.78	0.947	0.947	0.947	0.946	156	105	80	70
0.77	0.948	0.948	0.949	0.947	166	114	89	75
0.76	0.949	0.950	0.950	0.949	171	118	92	78
0.75	0.950	0.950	0.949	0.950	171	119	93	79

0.74	0.950	0.950	0.950	0.950	166	115	91	78
0.73	0.950	0.951	0.951	0.951	156	107	86	73
0.72	0.951	0.949	0.950	0.951	143	95	76	65
0.71	0.949	0.949	0.950	0.950	124	75	59	50

1
2 A second set of simulations was run to assess the performance of the method under the null
3 hypothesis where $\gamma_1 = \frac{TPR_A}{TPR_B} = 1$. Table 3 shows the cell probabilities for these simulations.
4 Rather than report across the entire range only the minimum, 50% and maximum levels of
5 *TPPR* are reported.

6 **Table 3** – Simulation settings to estimate Type I error

p_A	p_B	p_C	p_D	TPR_A	TPR_B	γ
0.81	0.045	0.045	0.10	0.855	0.855	1
0.76	0.095	0.095	0.05	0.855	0.855	1
0.71	0.145	0.145	0.00	0.855	0.855	1

7
8 Table 4 shows the type I error rate, mean sample size, bias, coverage and RMSE of simulated
9 sample sizes under various simulation settings. At all levels of conditional dependence and at
10 all interim sample sizes the type I error rate is close to the specified levels. Again, the
11 inference to be made from the RMSE value is that a larger sample size provides a more
12 accurate estimate of the full sample size required, reducing the extent to which an experiment
13 will be over or underpowered in individual cases. The bias and coverage also appear to be at
14 acceptable levels.

15 **Table 4** – Type I error rate, Mean sample size (S.D.), bias, coverage and RMSE of simulated
16 sample sizes under various simulation settings.

Type I error rate

	N=50	N=100	N=150	N=200
TPPR				
0.81	0.050	0.050	0.050	0.050
0.76	0.050	0.050	0.050	0.050
0.71	0.050	0.050	0.050	0.050

17

Mean sample size

Bias

N=50	N=100	N=150	N=200	N=50	N=100	N=150	N=200
------	-------	-------	-------	------	-------	-------	-------

TPPR									
0.81	304(121)	298(78)	297(61)	296(52)	0.00207	0.00224	0.00198	0.00186	
0.76	463(159)	457(107)	454(84)	453(71)	0.00147	0.00110	0.00100	0.00087	
0.71	627(118)	631(74)	630(58)	629(50)	0.00021	0.00023	0.00005	-0.00019	

1

TPPR	<u>Coverage</u>				<u>RMSE</u>			
	N=50	N=100	N=150	N=200	N=50	N=100	N=150	N=200
0.81	0.952	0.951	0.950	0.951	164	130	119	115
0.76	0.950	0.949	0.949	0.949	168	117	95	83
0.71	0.948	0.949	0.950	0.949	118	74	59	50

2

3 Table 5 gives the results of a range of simulations undertaken at various values of TPR_A ,
4 TPR_B and π in both true alternative and null cases. Regarding the best sample size to specify
5 at the interim, a possible balance to be struck between a suitably large interim sample, which
6 would increase the precision of the measure of conditional dependence, and minimising the
7 overall experimental sample size would be to take the minimal possible sample size for the
8 experiment as a whole at the interim. In this way, the interim sample could never be larger
9 than the overall required sample size, which means that it is impossible to collect more data
10 than actually needed. Yet, the minimum possible overall sample size represents a significant
11 proportion of the total experimental sample size. Thus, for the values in Table 5, the sample
12 size re-estimate was conducted at the number implied by equation 2, when $TPPR$ is at
13 maximal value given the marginals. The maximum positive, mid-range and maximum
14 negative levels of $TPPR$ were reported to show a range of values across different levels of
15 $TPPR$. The mean sample size is provided in parentheses in order to allow intuition about the
16 reduction in the sample size this method brings. In all cases, where data were generated under
17 the true alternative hypothesis, the simulated power is above or very close to the nominal
18 value. Furthermore, in all cases where data were generated under the true null hypothesis the
19 size is close to the nominal value. Comparing the mean sample sizes given for the maximal
20 and mid-point $TPPR$ s against the fixed values that would be used if Alonzo et al [19] had

1 been followed we can see that the sample size re-estimation method outlined above can
 2 dramatically reduce the required sample size to power an experiment to the minimum of a
 3 nominal level.

4
 5
 6
 7
 8

9 **Table 5** – Simulated type I and II error rates and fixed maximal sample size values under
 10 various true values of TPR_A , TPR_B and prevalence across various levels of conditional
 11 dependence (average sample size given in brackets)

	TPRb	TPRa	prev = 0.1			prev=0.3			prev=0.5		
			Alternative	Null	Fixed	Alternative	Null	Fixed	Alternative	Null	Fixed
Maximum positive TPR	0.5	0.6	0.979(871)	0.049(1255)	7084	0.977(289)	0.05(417)	2361	0.977(172)	0.048(249)	1417
	0.5	0.7	0.98(434)	0.047(616)	1585	0.98(143)	0.047(204)	528	0.978(85)	0.049(122)	317
	0.5	0.8	0.986(297)	0.048(401)	622	0.985(98)	0.047(133)	207	0.985(58)	0.048(79)	124
	0.5	0.9	0.991(232)	0.048(279)	303	0.99(76)	0.046(92)	101	0.989(45)	0.046(54)	61
	0.6	0.7	0.976(858)	0.047(1244)	5505	0.975(284)	0.05(414)	1835	0.975(170)	0.047(248)	1101
	0.6	0.8	0.978(431)	0.048(605)	1185	0.979(142)	0.049(200)	395	0.976(84)	0.048(119)	237
	0.6	0.9	0.984(297)	0.045(375)	442	0.984(98)	0.047(124)	147	0.985(58)	0.046(74)	88
	0.7	0.8	0.974(846)	0.048(1222)	3930	0.973(281)	0.05(410)	1310	0.971(167)	0.049(245)	786
	0.7	0.9	0.979(431)	0.049(575)	789	0.978(142)	0.046(190)	263	0.976(84)	0.049(114)	158
	0.8	0.9	0.971(837)	0.048(1195)	2357	0.971(277)	0.051(398)	786	0.97(165)	0.049(238)	471
50%TPPR	0.5	0.6	0.802(3974)	0.05(4050)	7084	0.806(1321)	0.049(1347)	2361	0.8(792)	0.049(807)	1417
	0.5	0.7	0.82(1013)	0.048(1082)	1585	0.822(336)	0.52(358)	528	0.818(200)	0.05(214)	317
	0.5	0.8	0.856(462)	0.049(517)	622	0.854(153)	0.046(171)	207	0.854(91)	0.048(102)	124
	0.5	0.9	0.911(270)	0.043(298)	303	0.908(89)	0.046(98)	101	0.905(52)	0.047(58)	61
	0.6	0.7	0.809(3175)	0.049(3277)	5505	0.805(1056)	0.05(1090)	1835	0.804(633)	0.049(653)	1101
	0.6	0.8	0.839(809)	0.051(891)	1185	0.838(268)	0.052(295)	395	0.835(160)	0.052(176)	237
	0.6	0.9	0.885(371)	0.046(416)	442	0.888(122)	0.047(137)	147	0.881(72)	0.047(82)	88
	0.7	0.8	0.809(2379)	0.053(2513)	3930	0.813(792)	0.052(836)	1310	0.812(474)	0.053(500)	786

	0.7	0.9	0.863(607)	0.05(687)	789	0.868(201)	0.052(228)	263	0.864(120)	0.051(136)	158
	0.8	0.9	0.832(1585)	0.05(1753)	2357	0.836(528)	0.051(583)	786	0.832(316)	0.051(349)	471
Maximal negative TPR	0.5	0.6	0.796(7105)	0.05(7109)	7084	0.797(2360)	0.49(2361)	2361	0.797(1413)	0.05(1414)	1417
	0.5	0.7	0.812(1608)	0.047(1609)	1585	0.804(533)	0.051(534)	528	0.81(318)	0.05(318)	317
	0.5	0.8	0.827(639)	0.05(641)	622	0.829(211)	0.049(212)	207	0.832(126)	0.049(126)	124
	0.5	0.9	0.87(312)	0.05(317)	303	0.867(103)	0.048(104)	101	0.868(61)	0.048(62)	61
	0.6	0.7	0.798(5522)	0.05(5519)	5505	0.796(1836)	0.049(1836)	1835	0.798(1099)	0.05(1099)	1101
	0.6	0.8	0.812(1204)	0.051(1205)	1185	0.818(399)	0.051(400)	395	0.812(238)	0.051(238)	237
	0.6	0.9	0.844(451)	0.048(456)	442	0.842(149)	0.047(151)	147	0.839(89)	0.047(90)	88
	0.7	0.8	0.799(3944)	0.049(3943)	3930	0.806(1311)	0.049(1311)	1310	0.803(785)	0.048(785)	786
	0.7	0.9	0.824(799)	0.052(803)	789	0.826(265)	0.05(266)	263	0.824(158)	0.05(159)	158
	0.8	0.9	0.806(2366)	0.052(2368)	2357	0.808(787)	0.05(788)	786	0.808(471)	0.051(471)	471

1
2

3 Application Revisited

4 Given the robustness of the proposed method of sample size recalculation described and
5 validated in simulation above, we return to apply it to the application presented earlier in this
6 paper. The cell probability values at maximum positive conditional dependence for diseased
7 patients under the specified values of TPR_A and TPR_B are $\hat{p}_1 = 0.81, \hat{p}_2 = 0.09, \hat{p}_3 = 0,$
8 $\hat{p}_4 = 0.1$. The cell probability values at maximum negative conditional dependence for
9 diseased patients under the specified values of TPR_A and TPR_B are $\hat{p}_1 = 0.71, \hat{p}_2 = 0.19,$
10 $\hat{p}_3 = 0.10, \hat{p}_4 = 0$. Table 7 shows an example range of the permissible values under the
11 specified values of TPR_A and TPR_B . Given this, we can create a likelihood of our observed
12 interim data having come from each possible configuration of the alternative hypothesis using
13 equation 3.

14 **Table 7** – Example range of cell probabilities based on: $TPR_A = 0.9$ and $TPR_B = 0.81$

p_1	p_2	p_3	p_4	TPR_A	TPR_B
0.81	0.09	0.00	0.10	0.9	0.81
0.80	0.10	0.01	0.09	0.9	0.81
...
0.72	0.18	0.09	0.01	0.9	0.81
0.71	0.19	0.10	0.00	0.9	0.81

1
2 Applying the method outlined in section 3, we take; $\widehat{TPR}_A = 0.9$, $\widehat{TPR}_B = 0.81$, observed
3 $\hat{n}_A = 66$, $\hat{n}_B = 3$, $\hat{n}_C = 3$, $\hat{n}_D = 10$ and $\hat{n}_E + \hat{n}_F + \hat{n}_G + \hat{n}_H = 105$, implying $\hat{\pi} =$
4 0.439 . Using equation 4 the maximum likelihood value of \widehat{TPPR} is 0.793 . Given the fact
5 that π is binomially distributed, the maximum likelihood estimate for the prevalence is equal
6 to the observed prevalence, $\hat{\pi}$. Taking these values and inserting them into equation 1 we get
7 the value for the sample size required for sensitivity as **275**. Taking $\widehat{TNRA} = 0.8$ and
8 $\widehat{TNRB} = 0.66$, with the observed values $\hat{n}_E = 21$, $\hat{n}_F = 4$, $\hat{n}_G = 11$, $\hat{n}_G = 69$ and $\hat{n}_A +$
9 $\hat{n}_B + \hat{n}_C + \hat{n}_D = 82$, implying $1 - \hat{\pi} = 0.561$. Using equation 3 to derive the maximum
10 likelihood of the cell probabilities for specificity we estimate that $\widehat{TNNR} = 0.635$. Inserting
11 these values into equation 2 gives us a sample size estimate of **136**. Thus, the updated sample
12 size, in order to use the interim information about the conditional dependence between the
13 tests and to preserve a minimal nominal power of 0.8 should be **275**.

14 **Discussion**

15 This paper has presented a robust method of sample size re-estimation for use in paired
16 diagnostic accuracy studies where the conditional independence between the two tests may be
17 unknown or inaccurately estimated at the start of the study. In terms of the recommendation
18 of sample size estimation for the experiment as a whole a specific protocol is suggested given
19 the results. Rather than basing the estimate for the experiment as a whole on the case where
20 there is the maximal negative conditional dependence between tests – thus the largest
21 possible sample size - as suggested in Alonzo et al [19], we would suggest an alternative
22 strategy, the robustness of which is highlighted in Table 5. Specifically, initially estimating
23 the sample size at the maximal positive conditional dependence between tests, i.e. using
24 $TPPR = TPR_B$ - giving the smallest possible sample size - then, re-estimating the final
25 sample size using the method simulated in Table 5. As long as the initial estimate for

1 prevalence is close to accurate, this protocol is deemed appropriate as it balances the risk of
2 collecting more participants than might actually be needed with collecting the most
3 information about the true conditional dependence at the interim. Table 5 provides strong
4 evidence for the integrity of this method in providing at minimum the nominal power while
5 reducing the sample size when we have a higher than maximally negative true conditional
6 dependence. Should the interim sample size be some other value, the maximum likelihood
7 method will still be appropriate, although it should be kept in mind that the larger the interim
8 sample size, as a proportion of the total possible sample size, the more accurate the interim
9 sample size estimates will be, for individual cases.

10 Interestingly, the sample size values in the table seem to be somewhat greater, even when
11 using our method than those typically seen in the literature in diagnostic test accuracy studies,
12 see for example van Enst *et al.* [27] Although it is difficult to know the specifics of the 859
13 studies mentioned in the van Enst collection of meta-analyses, e.g. clinically significant
14 differences, sample size estimation and hypothesis testing procedures, it is striking that the
15 50% covariance sample size is only 87 (IQR 45-185) participants. Very few of our sample
16 sizes in Table 5 are this low for the size of effect (ratios) we are considering, even using our
17 method of sample size reduction. It may be that many diagnostic accuracy studies
18 commissioned do not carefully consider their sample sizes.

19 While the method discussed here of estimating the conditional dependence between the tests
20 via maximum likelihood, given constraints imposed by the specified marginals and under a
21 multinomial model, is pertinent to paired diagnostic accuracy tests, there is little reason why
22 similar processes could not be extended to similar problems. The kernel of the method,
23 maximum likelihood estimation of the parameter related to the conditional dependence using

1 a constrained multinomial model, is equally valid in other applications involving sample size
2 re-estimation for paired binary 2x2 tables.

3 **Conclusions**

4 In this paper we have described a sample size re-estimation procedure that can be applied in
5 an interim analysis for a diagnostic test study that is comparing two methods of testing on
6 patients that are being followed up over a period of time. The procedure uses information on
7 the levels of conditional dependence between the two tests at the interim in order to refine the
8 required sample size for a paired diagnostic accuracy study with a binary response. Evidence
9 from simulations has been provided to demonstrate its functionality under various parameter
10 values thought to reflect a range of commonly occurring situations. The procedure can be
11 applied in the case of paired comparative diagnostic accuracy studies in order to more
12 accurately gauge the sample size required for a given power thereby reducing both the costs
13 associated with this kind of study and also the burden on patients.

14

15

16

17

18

19

20

1

2

3

4

5

6

7

8 **List of abbreviations**

9 CT: computed tomography

10 PET: Positron Emission Tomography

11 RMSE: Root Mean Squared Error

12 TPR_A : True positive rate for test A

13 TPR_B : True positive rate for test B

14 $TPPR$: True positive positive rate

15 TNR_A : True negative rate for test A

16 TNR_B : True negative rate for test B

17 $TNNR$: True negative negative rate

18

1

2

3

4

5

6

7

8

9 **Declarations**

10 **1. Ethics approval and consent to participate**

11 REC reference: 10 H1017 8. The original application was made to North West 1
12 Research Ethics Committee - Cheshire. This committee has subsequently been
13 superseded by the North West - Greater Manchester East Research Ethics Committee.
14 Consent to participate was provided by all participants.

15 **2. Consent for publication**

16 Not applicable

17 **3. Availability of data and material**

18 All data used to illustrate the method can be found in Table 6 of this paper

19 **4. Competing interests**

1 The authors declare that they have no competing interests.

2 **5. Funding**

3 This project was funded by the NIHR Health Technology Assessment Programme, REF:
4 08/29/02.

5 **6. Authors' contributions**

6 GL conceived the study. The method was devised by GM, overseen by GL and AT. The
7 manuscript was written by GM. All authors contributed to the reviewing and revising of
8 the manuscript. All authors read and approved the final manuscript.

9 **7. Acknowledgements**

10 Not applicable

11

12

13

14

15

16

17

18

19

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21

References

1. Knottnerus JA, Weel C Van, Muris JWM. Evaluation of diagnostic procedures. *Br. Med. J.* 2002;324:477–80.
2. Swets JA, Pickett RM. *Evaluation of Diagnostic Systems*. New York: Academic Press; 1982.
3. Zhou XH, Obuchowski NA, McClish DK. *Statistical Methods in Diagnostic Medicine*. New York: Wiley; 2002.
4. Freedman LS. Evaluating and comparing imaging techniques: a review and classification of study designs. *Br. J. Radiol.* 1987;60:1071–81.
5. Gerke O, Vach W, Høilund-Carlsen PF. PET/CT in Cancer. *Methods Inf. Med.* [Internet]. 2008 [cited 2014 Oct 15];470–9. Available from: <http://www.schattauer.de/index.php?id=1214&doi=10.3414/ME0540>
6. Gould AL. Sample size re-estimation: recent developments and practical considerations.

- 1 Stat. Med. 2001;20:2625–43.
- 2 7. Wittes J, Brittain E. The role of internal pilot studies in increasing the efficiency of clinical
3 trials. Stat. Med. 1990;9:65–72.
- 4 8. Shih WJ. Sample size reestimation in clinical trials. In: Peace K, editor. Biopharm. Seq.
5 Stat. Appl. New York: Marcel Dekker; 1992. p. 285–301.
- 6 9. Shih WJ. Sample size reestimation for triple blind clinical trials. Drug Inf. J. 1993;27:761–
7 4.
- 8 10. Gould AL, Shih WJ. Sample size re-estimation without unblinding for normally
9 distributed outcomes with unknown variance. Commun. Stat. Methods. 1992;21:2833–53.
- 10 11. Birkett MA, Day SJ. Internal Pilot Studies for Estimating Sample Size. Stat. Med.
11 1994;13:2455–63.
- 12 12. Gould AL. Interim Analyses for Monitoring Clinical Trails that do not Affect Type I
13 Error Rates. Stat. Med. 1992;11:55–66.
- 14 13. Herson J, Wittes J. The Use of Interim Analysis for Sample Size Adjustment. Drug Inf. J.
15 1993;27:761–4.
- 16 14. Shih WJ, Zhao P. Design for Sample Size Re-estimation with Interim Data for Double-
17 Blind Clinical Trails. Stat. Med. 1997;16:1913–23.
- 18 15. Proschan MA. Two-stage sample size re-estimation based on a nuisance parameter: a
19 review. J. Biopharm. Stat. 2005;15:559–74.
- 20 16. Gerke O, Høilund-carlsen PF, Poulsen MH, Vach W. Interim analyses in diagnostic
21 versus treatment studies : differences and similarities. Am. J. Nucl. Med. Mol. Imaging.
22 2012;2:344–52.
- 23 17. Newcombe RG. Improved Confidence Intervals for the Difference between Binomial

- 1 Proportions Based on Paired Data. *Stat. Med.* 1998;17:2635–50.
- 2 18. Tango T. Equivalence Test and Confidence Interval for the Difference in the Porportions
3 Based on Paired Data. *Stat. Med.* 1998;17:891–908.
- 4 19. Alonzo TA, Pepe MS, Moskowitz CS. Sample Size Calculations for Comparative Studies
5 of Medical Tests for Detecting Presence of Disease. *Stat. Med.* 2002;21:835–52.
- 6 20. Lu Y, Jin H, Genant HK. On the Non-Inferiority of a Diagnostic Test Based on Paired
7 Observations. *Stat. Med.* 2006;3:227–79.
- 8 21. Moskowitz CS, Pepe MS. Comparing the Predictive Values of Diagnostic Tests: Sample
9 Size and Analysis for PAired Study Designs. *Clin. trials.* 2006;3:272–9.
- 10 22. Bonett DG, Price RM. Confidence Intervals for a Ratio of Binomial Proportions Based on
11 Paired Data. *Stat. Med.* 2006;25:3039–47.
- 12 23. Vacek P. The effect of conditional dependence on the evaluation of diagnostic tests.
13 *Biometrics.* 1985;41:959–68.
- 14 24. van Smeden M, Naaktgeboren CA, Reitsma JB, Moons KGM, Groot JAH De. Latent
15 Class Models in Diagnostic Studies When There is No Reference Standard — A Systematic
16 Review. *Am. J. Epidemiol.* 2014;179:423–31.
- 17 25. Schiller I, Smeden M Van, Hadgu A, Libman M, Reitsma B, Dendukuri N. Bias due to
18 composite reference standards in diagnostic accuracy studies. *Stat. Med.* 2016;35:1454–70.
- 19 26. Royston P. Exact conditional and unconditional sample size for pair-matched studies with
20 binary outcome: a practical guide. *Stat. Med.* 1993;12:699–712.
- 21 27. van Enst WA, Naaktgeboren CA, Ochodo EA, Groot JAH De, Leeftang MM, Reitsma
22 JB, et al. Small-study effects and time trends in diagnostic test accuracy meta-analyses : a
23 meta-epidemiological study. *Syst. Rev.* 2015;4.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35

Appendix – R code for maximum likelihood sample size re-estimation.

```
ss.est.mle <- function(obs.a, obs.b, obs.c, obs.d, obs.x,  
tpra, tprb, alpha, beta){  
mle.tppr <- function(theta.1, obs.a, obs.b, obs.c, obs.d,  
tpra, tprb)  
{-((obs.a*log(theta.1)) + obs.b*log(tpra-theta.1) +  
obs.c*log(tprb-theta.1) + obs.d*log(-tpra-tprb+theta.1+1))}  
tppr <- optim(par=(tpra+(tprb-(1-tpra)))/2 ,fn=mle.tppr,  
obs.a=obs.a, obs.b=obs.b, obs.c=obs.c, obs.d=obs.d, tpra=tpra,  
tprb=tprb, method = "Brent", lower=(tprb-(1-tpra)),  
upper=tprb)$par  
obs.prev <-  
(obs.a+obs.b+obs.c+obs.d)/(obs.a+obs.b+obs.c+obs.d+obs.x)  
alonzo <- function(lambda, prev ,beta, alpha, tprb,  
gam1){(((qnorm(1-beta) + qnorm(1-alpha))/log(gam1))^2 *  
((gam1+1) * tprb)-(2 * lambda))/(gam1*tprb^2))/prev}  
gam1 <- tpra/tprb  
ss.est <- alonzo(tppr, obs.prev, beta = beta, alpha = alpha,  
tprb = tprb, gam1 = gam1 )
```

```
1 return(ss.est)}
2
3
4 ### Example sensitivity
5
6 ss.est.mle(obs.a=66, obs.b=3, obs.c=3, obs.d=10, obs.x=105,
7 tpra=0.9, tprb=0.81, alpha=0.025, beta=0.2)
8
9 ### Example specificity
10
11 ss.est.mle(obs.a=69, obs.b=11, obs.c=4, obs.d=21, obs.x=82,
12 tpra=0.8, tprb=0.66, alpha=0.025, beta=0.2)
13
```