

# Optimality of group testing in the presence of misclassification

BY AIYI LIU

*Biostatistics and Bioinformatics Branch, Eunice Kennedy Shriver National Institute of Child Health and Human Development, Rockville, Maryland 20852, U.S.A.*

liua@mail.nih.gov

CHUNLING LIU

*Department of Applied Mathematics, Hong Kong Polytechnic University, Hong Kong, China*

Catherine.Chunling.Liu@inet.polyu.edu.hk

ZHIWEI ZHANG AND PAUL S. ALBERT

*Biostatistics and Bioinformatics Branch, Eunice Kennedy Shriver National Institute of Child Health and Human Development, Rockville, Maryland 20852, U.S.A.*

zhangz7@mail.nih.gov albertp@mail.nih.gov

## SUMMARY

Several optimality properties of Dorfman's (1943) group testing procedure are derived for estimation of the prevalence of a rare disease whose status is classified with error. Exact ranges of disease prevalence are obtained for which group testing provides more efficient estimation when group size increases.

*Some key words:* Binary outcome; Maximum likelihood estimation; Pooling; Prevalence; Sensitivity; Specificity.

## 1. INTRODUCTION

Aiming at more efficient screening of a rare disease, Dorfman (1943) proposed to test for the syphilis antigen by first testing pooled blood samples, followed by retesting individuals in groups found to be infected. This strategy and its variations developed later, often referred to as group testing or pooled testing, have received substantial attention for efficient identification of an event or estimation of the probability that the event occurs; see Sobel & Groll (1959), Sobel & Elashoff (1975), Le (1981), Gastwirth & Hammick (1989), Chen & Swallow (1990), Farrington (1992), Gastwirth & Johnson (1994), Hughes-Oliver & Swallow (1994), Litvak et al. (1994), Tu et al. (1995), Barcellos et al. (1997), Brookmeyer (1999), Hung & Swallow (1999), Hughes-Oliver & Rosenberger (2000) and Tebbs & Swallow (2003). An attractive feature of group testing is that retesting on individuals is not necessary if one is only interested in estimation of the probability of a positive test.

Most of the developments and applications of group testing have assumed that the disease status can be accurately determined without error. In this case, use of group testing can substantially reduce the cost but will always yield estimation of prevalence less efficient than that based on the fully observed data, in which the disease status is determined for each individual. A number of authors have investigated group testing strategies with misclassification of the disease status; see, among others, Graff & Roeloffs (1972), Hwang (1976), Burns & Mauro (1987) and Xie et al. (2001). Tu et al. (1995) showed that when the total number of subjects is fixed, the asymptotic variance function of the resulting estimator decreases as the number of subjects in each group increases, assuming that the prevalence of the disease is small enough.

It is not clear, however, exactly how small the prevalence needs to be for group testing to become more efficient, since [Tu et al. \(1995\)](#) only considered the limiting behaviour of the variance.

In this paper we further investigate the optimality properties of the group testing strategy in estimating the prevalence of a disease. We show that, when the disease status is measured with error, group testing with moderate group sizes provides more efficient estimation than the fully observed data over a wide range of disease prevalences. When the number of groups is fixed, group testing also prevails over the one-subject-per-group random sampling design for moderate disease prevalence.

## 2. OPTIMALITY OF THE ESTIMATION

### 2.1. Maximum likelihood estimation

We assume that the status of a disease  $D$  is measured with error, with specificity  $\pi_0 = \text{pr}(M = 0 \mid D = 0)$  and sensitivity  $\pi_1 = \text{pr}(M = 1 \mid D = 1)$ , where  $M = 0, 1$  is the observed status of  $D$ . For the classification to be of practical use we assume that  $1/2 < \pi_0, \pi_1 \leq 1$ .

Suppose that  $k$  samples are pooled and that instead of observing  $M_1, \dots, M_k$ , we observe  $\tilde{M} = \max(M_1, \dots, M_k)$ . Thus, if  $\tilde{M} = 0$ , then  $M_i = 0$  for each  $i = 1, \dots, k$ . If  $\tilde{M} = 1$ , then we conclude that  $M_i = 1$  for at least one subject in the group. We assume that misclassification of the disease is nondifferentiable, independent of pooling and group sizes. This is a common assumption in practice. Write  $D_i = 1$  if the  $i$ th subject has the disease and 0 if otherwise. Define  $p = \text{pr}(D_1 = 1)$ , the prevalence of the disease, and assume that  $0 < p < 1$ . It then follows from [Tu et al. \(1995\)](#) that  $\text{pr}(\tilde{M} = 1) = \pi_1 - r(1 - p)^k \in (1 - \pi_0, \pi_1)$ , where  $r = \pi_0 + \pi_1 - 1$ .

Let  $n$  be the number of groups being tested, each of size  $k$ , the number of subjects in the group. Let  $n_1 \leq n$  be the number of groups that test positive and define  $\lambda = n_1/n$ , the observed proportion of positive groups. Then the maximum likelihood estimate is

$$\hat{p} = 1 - \left[ \frac{\pi_1 - \min\{\pi_1, \max(1 - \pi_0, \lambda)\}}{r} \right]^{1/k} \quad (1)$$

with asymptotic variance

$$\text{var}(\hat{p}) = \sigma^2(p, k, n) = \frac{\{\pi_1 - r(1 - p)^k\}\{r(1 - p)^k + 1 - \pi_1\}}{nr^2k^2(1 - p)^{2(k-1)}}. \quad (2)$$

See [Tu et al. \(1995\)](#) and [Hepworth \(1996\)](#), among others.

From the variance formula (2) we notice that for fixed group size  $k$  and misclassification rates,  $\text{var}(\hat{p})$  decreases as the number  $n$  of groups increases. Therefore the precision of  $\hat{p}$  can be improved by increasing the number of groups, which may be infeasible in epidemiological studies due to cost constraints.

### 2.2. Optimality with a fixed number of groups

Of particular interest is the relative efficiency of group testing to a random sample of size  $n$ , or the one-sample-per-group design, assuming that the number  $n$  of groups is fixed. This situation often occurs in practice when only a limited number of test-kits or assays are available due to cost constraints. The random sampling design randomly selects  $n$  subjects from the population, with disease status observed for each subject. Thus the random sampling design results in an estimate of  $p$ , given by (1) with  $k = 1$ , with variance  $\sigma^2(p, 1, n) = (\pi_0 - rp)(1 + rp - \pi_0)/(nr^2)$ , derived from (2) by straightforward algebraic manipulation.

**PROPOSITION 1.** Consider  $\sigma^2(p, k, n)$  as defined in (2). Let  $n, k_1$  and  $k_2$  be integers and  $1 \leq k_2 < k_1$ . Then  $\sigma^2(p, k_1, n)/\sigma^2(p, k_2, n)$  is a monotone increasing function of  $p$ .

*Proof.* The derivative of  $\log\{\sigma^2(p, k_1, n)/\sigma^2(p, k_2, n)\}$  with respect to  $p$  is given by  $f(k_1; 1 - \pi_1) + f(k_1; -\pi_1)$ , where

$$f(x; s) = \frac{s}{1-p} \left\{ \frac{x}{r(1-p)^x + s} - \frac{k_2}{r(1-p)^{k_2} + s} \right\} \quad (x \geq k_2).$$

Since the derivative of  $f(x; 1 - \pi_1)$  with respect to  $x$  is positive,  $f(x; 1 - \pi_1)$  is monotone increasing in  $x$  and thus  $f(k_1; 1 - \pi_1) > f(k_2; 1 - \pi_1) = 0$ . Furthermore, the derivative of  $f(x; -\pi_1)$  with respect to  $x$  has the same sign as  $g(p) = \pi_1 - r(1-p)^x + rx(1-p)^x \log(1-p)$ , whose derivative with respect to  $p$  is  $-rx^2(1-p)^{x-1} \log(1-p) > 0$ . Thus  $g(p)$  is monotone increasing, implying  $\lim_{p \rightarrow 1} g(p) = \pi_1 > g(p) > \lim_{p \rightarrow 0} g(p) = 1 - \pi_0 \geq 0$ . Therefore  $f(x; -\pi_1)$  is also monotone increasing in  $x$  and thus  $f(x; -\pi_1) > f(k_2; -\pi_1) = 0$ . This completes the proof.  $\square$

**PROPOSITION 2.** *Assume that the conditions in Proposition 1 hold. Then there exists some unique  $0 < p_0 < 1$  such that  $\sigma^2(p_0, k_1, n) = \sigma^2(p_0, k_2, n)$  and*

$$\sigma^2(p, k_1, n) \begin{cases} < \sigma^2(p, k_2, n) & (p < p_0), \\ > \sigma^2(p, k_2, n) & (p > p_0). \end{cases}$$

*Proof.* Straightforward manipulations lead to

$$\lim_p \frac{\sigma^2(p, k_1, n)}{\sigma^2(p, k_2, n)} = \begin{cases} \infty & (p \rightarrow 1), \\ k_2^2/k_1^2 < 1 & (p \rightarrow 0, \pi_0 < 1), \\ k_2/k_1 < 1 & (p \rightarrow 0, \pi_0 = 1). \end{cases}$$

The proposition thus follows because  $\sigma^2(p, k_1, n)/\sigma^2(p, k_2, n)$  is a continuous function of  $p$  and is monotone increasing in  $p$ .  $\square$

With respect to the random sampling we have the following result.

**THEOREM 1.** *Let  $k \geq 2$  and  $0 < p_0 < 1$  be such that  $\sigma^2(p_0, k, n) = \sigma^2(p_0, 1, n)$ . Then group testing with group size  $k$  and number  $n$  of groups is more efficient in estimating the disease prevalence  $p$  than a random sample of size  $n$  if and only if  $p < p_0$ .*

### 2.3. Optimality with a fixed number of subjects

Next we consider the situation in which the total number of subjects  $nk$  is fixed. The issue under investigation is then to find  $k$  for which group testing improves the precision in estimating  $p$  when compared with the fully observed data, that is, data with disease status observed for each of the  $nk$  subjects. The problem reduces to comparing  $\sigma^2(p, k, n)$  with  $\sigma^2(p, 1, nk)$ . Theorem 3 of [Tu et al. \(1995\)](#) proved that for  $\pi_0 < 1$  and fixed  $nk$  the variance function (2) is monotone decreasing in  $k$  when  $p$  is small. It follows that for  $\pi_0 < 1$  and small  $p$ ,  $\sigma^2(p, k, n) < \sigma^2(p, 1, nk)$  for  $k \geq 2$ . However, the authors did not discuss how small  $p$  needs to be for the inequality to hold. We provide the following results to address this.

**PROPOSITION 3.** *Consider  $\sigma^2(p, k, n)$  as defined in (2) where  $0 < p < 1$ . Let  $k_1, k_2, n_1$  and  $n_2$  be integers such that  $k_1 > k_2 \geq 1$  and  $n_1 k_1 = n_2 k_2$ .*

- (a) *If  $\pi_0 = 1$ , then  $\sigma^2(p, k_1, n_1) > \sigma^2(p, k_2, n_2)$ .*
- (b) *If  $\pi_0 < 1$ , then there exists some unique  $0 < p_0 < 1$  such that  $\sigma^2(p_0, k_1, n_1) = \sigma^2(p_0, k_2, n_2)$  and*

$$\sigma^2(p, k_1, n_1) \begin{cases} < \sigma^2(p, k_2, n_2) & (p < p_0), \\ > \sigma^2(p, k_2, n_2) & (p > p_0). \end{cases}$$

Table 1. Values of  $p_0$  below which group testing is more efficient than random sampling/fully observed data

		$k = 2$			
		$\pi_1 = 0.85$	0.90	0.95	0.99
$\pi_0 = 0.85$		0.543/0.247	0.568/0.255	0.606/0.265	0.662/0.276
0.90		0.537/0.216	0.563/0.222	0.602/0.229	0.657/0.238
0.95		0.531/0.167	0.557/0.171	0.596/0.176	0.651/0.181
0.99		0.524/0.085	0.552/0.086	0.591/0.088	0.646/0.090
		$k = 5$			
		$\pi_1 = 0.85$	0.90	0.95	0.99
$\pi_0 = 0.85$		0.357/0.156	0.376/0.161	0.405/0.169	0.457/0.177
0.90		0.352/0.136	0.371/0.140	0.401/0.146	0.451/0.152
0.95		0.345/0.106	0.365/0.108	0.395/0.112	0.445/0.116
0.99		0.338/0.054	0.359/0.055	0.390/0.056	0.439/0.057
		$k = 10$			
		$\pi_1 = 0.85$	0.90	0.95	0.99
$\pi_0 = 0.85$		0.239/0.105	0.252/0.109	0.272/0.115	0.312/0.122
0.90		0.235/0.092	0.248/0.095	0.268/0.100	0.307/0.105
0.95		0.229/0.072	0.242/0.074	0.263/0.077	0.302/0.080
0.99		0.223/0.037	0.236/0.038	0.258/0.039	0.296/0.040

*Proof.* If  $n_1k_1 = n_2k_2$ , then

$$\frac{\sigma^2(p, k_1, n_1)}{\sigma^2(p, k_2, n_2)} = \frac{k_1 \sigma^2(p, k_1, n_1)}{k_2 \sigma^2(p, k_2, n_1)}$$

The variance ratios do not involve  $n_1$  and  $n_2$ , so it follows from Proposition 1 that  $\sigma^2(p, k_1, n_1)/\sigma^2(p, k_2, n_2)$  is also a monotone increasing function in  $p$ . Moreover

$$\lim_p \frac{\sigma^2(p, k_1, n_1)}{\sigma^2(p, k_2, n_2)} = \begin{cases} \infty & (p \rightarrow 1), \\ k_2/k_1 < 1 & (p \rightarrow 0, \pi_0 < 1), \\ 1 & (p \rightarrow 0, \pi_0 = 1). \end{cases}$$

Thus (a) and (b) follow. □

**THEOREM 2.** Let  $n$  and  $k$  be integers and  $0 < p_0 < 1$  be such that  $\sigma^2(p_0, k, n) = \sigma^2(p_0, 1, nk)$ . Then group testing with size  $k$  and number  $n$  of groups is more efficient in estimating the disease prevalence  $p$  than a random sample of size  $nk$  if and only if  $p < p_0$  and  $\pi_0 < 1$ .

Therefore, when the disease status is measured with error, group testing yields more efficient estimation of the disease prevalence than both random sampling and fully observed data, if the disease prevalence is relatively low.

The values of  $p_0$  in Theorems 1 and 2 depend on the group size  $k$ , the sensitivity  $\pi_1$  and specificity  $\pi_0$ , but not on the number  $n$  of groups. For selected values of sensitivity  $\pi_1$ , specificity  $\pi_0$  and group size  $k$ , values of  $p_0$  in Theorem 1 and Theorem 2 are presented in Table 1. Both cases reveal some common features. For fixed  $k$  and  $\pi_1$ ,  $p_0$  decreases as  $\pi_0$  increases; for fixed  $k$  and  $\pi_0$ ,  $p_0$  increases as  $\pi_1$  increases; and for fixed  $\pi_0$  and  $\pi_1$ ,  $p_0$  decreases as  $k$  increases. The range of  $p_0$  values for the group testing to be more efficient is wider if compared with a random sample of size  $n$  than if compared with fully observed data. This reflects the fact that for fixed group sizes and misclassification rates, the precision of estimation increases as the number of subjects increases.

### 3. EXAMPLE: ESTIMATION OF HUMAN IMMUNODEFICIENCY VIRUS PREVALENCE

Human immunodeficiency virus is a lentivirus that causes acquired immunodeficiency syndrome, and is present as both free virus particles and as a virus within infected immune cells. According to the Morbidity and Mortality Weekly Report of the [United States Centers for Disease Control and Prevention \(2009\)](#), the prevalence per 100 000 population of human immunodeficiency virus infection in the U.S.A. is about 2388.2 among black men and 394.6 among white men, yielding  $p \approx 0.024$  and  $0.004$ , respectively.

Suppose we wish to estimate the prevalence by applying group testing to a given number of black men and white men, and that the presence of the human immunodeficiency virus is detected by an enzyme-linked immunosorbent assay, with 95% sensitivity and 90% specificity, close to what was reported in [Weiss et al. \(1985\)](#). With  $p = 0.024$  or  $0.004$ , group testing with group sizes up to  $k = 15$  yields more efficient estimation of the prevalence than fully observed data by assaying each individual.

However, the group sizes for which group testing is more efficient differ between the two groups. By plotting the variance (2) against group size, we found that group testing achieves minimum variance with group size  $k = 15$  for black men and  $k = 93$  for white men; see Fig. 1 in the Supplementary Material. Thus a group testing strategy that is optimal for estimating prevalence among white men may suffer from loss of efficiency if used to estimate prevalence among black men. For example, the variance with  $k = 90$  is more than three times the variance with  $k = 15$ ; see Fig. 2 in the Supplementary Material.

Another example provided in the Supplementary Material concerns estimation of gene-environment interaction in a case-control study.

## 4. DISCUSSION

### 4.1. Finite-sample results

For a given sample size, exact variance, bias and mean squared error of the maximum likelihood estimate can be calculated using binomial distributions. We found it extremely difficult to derive conditions under which one group testing strategy performs better than another in terms of the exact results. Furthermore, unlike the asymptotic case, we suspect that such a condition also involves the number of groups.

We further investigated the accuracy of the values of  $p_0$  for small sample sizes in terms of the exact variance, bias and mean squared error. The numerical results are presented in the Supplementary Material. They indicate that caution needs to be exercised when using the asymptotic results for relatively small sample sizes. It appears that the ranges of  $p$  over which group testing performs better are smaller than those based on large-sample variance.

### 4.2. Adaptive designs

The optimal group size depends on the disease prevalence  $p$  to be estimated and the sensitivity  $\pi_1$  and specificity  $\pi_0$ , which are often unknown. The values of these parameters must be specified in order to determine a proper group size  $k$ . However, the group size thus determined can result in much less efficient estimation if the parameters are misspecified. It is therefore desirable to develop procedures that are less sensitive to the specified values of the unknown parameters. To this end, adaptation of Stein's (1945) two-stage sampling procedure may be promising. [Hughes-Oliver & Swallow \(1994\)](#) and [Ridout \(1995\)](#) investigated the performance of such adaptive strategies for group testing without misclassification, and [Hepworth & Watson \(2009\)](#) presented methods to reduce the bias of the maximum likelihood estimate upon completion of such a scheme. For group testing with misclassification, we expect that these types of adaptive strategies may provide for efficient estimation when little a priori information exists about  $\pi_0$ ,  $\pi_1$  and  $p$ . This is an area of future research.

### 4.3. Differentiable misclassification

Throughout we have assumed that misclassification is known, and is unaffected by pool size, disease status, and other factors. These assumptions may well be violated in practice, with testing errors depending

on disease status and/or pool size; see, e.g., Zhang et al. (2008) and Cahoon-Young et al. (1989). Further research is needed to extend the results derived in the present paper to differentiable misclassification.

#### ACKNOWLEDGEMENT

This research was supported by the Intramural Research Program of the Eunice Kennedy Shriver National Institute of Child Health and Human Development of the National Institutes of Health. The authors thank the editor, associate editor and two referees for their constructive comments, and Dr Yaakov Malinovsky for helpful discussions.

#### SUPPLEMENTARY MATERIAL

Supplementary material available at *Biometrika* online includes more details on estimation of human immunodeficiency virus prevalence, an example on using group testing to estimate gene-environment interaction in case-control studies, and some finite-sample results.

#### REFERENCES

- BARCELLOS, L. F., KLITZ, W., FIELD, L. L., TOBIAS, R., BOWCOCK, A. M., WILSON, R., NELSON, M. P., NAGATOMI, J. & THOMSON, G. (1997). Association mapping of disease loci, by use of a pooled DNA genomic screen. *Am. J. Hum. Genet.* **61**, 734–47.
- BROOKMEYER, R. (1999). Analysis of multistage pooling studies of biological specimens for estimating disease incidence and prevalence. *Biometrics* **55**, 608–12.
- BURNS, K. C. & MAURO, C. A. (1987). Group testing with test error as a function of concentration. *Commun. Statist. A* **16**, 2821–37.
- CAHOON-YOUNG, B., CHANDLER, A., LIVERMORE, T., GAUDINO, J. & BENJAMIN, R. (1989). Sensitivity and specificity of pooled versus individual sera in a human immunodeficiency virus (mY) antibody prevalence study. *J. Clin. Microbiol.* **27**, 1893–5.
- CHEN, C. L. & SWALLOW, W. H. (1990). Using group testing to estimate a proportion, and to test the binomial model. *Biometrics* **46**, 1035–46.
- DORFMAN, R. (1943). The detection of defective members of large populations. *Ann. Math. Statist.* **14**, 436–40.
- FARRINGTON, C. P. (1992). Estimating prevalence by group testing using generalized linear models. *Statist. Med.* **11**, 1591–7.
- GASTWIRTH, J. L. & HAMMICK, P. A. (1989). Estimation of prevalence of a rare disease, preserving anonymity of subjects by group testing: application to estimating the prevalence of AIDS antibodies in blood donors. *J. Statist. Plan. Infer.* **22**, 15–27.
- GASTWIRTH, J. & JOHNSON, W. (1994). Screening with cost-effective quality control: potential applications to HIV and drug testing. *J. Am. Statist. Assoc.* **89**, 972–81.
- GRAFF, L. E. & ROELOFFS, R. (1972). Group testing in the presence of test errors: an extension of the Dorfman procedure. *Technometrics* **14**, 113–22.
- HEPWORTH, G. (1996). Exact confidence intervals for proportions estimated by group testing. *Biometrics* **52**, 1134–46.
- HEPWORTH, G. & WATSON, R. (2009). Debiased estimation of proportions in group testing. *J. R. Statist. Soc. C* **58**, 105–21.
- HUGHES-OLIVER, J. M. & SWALLOW, W. H. (1994). A two-stage adaptive group-testing procedure for estimating small proportions. *J. Am. Statist. Assoc.* **89**, 982–93.
- HUGHES-OLIVER, J. M. & ROSENBERGER, W. F. (2000). Efficient estimation of the prevalence of multiple rare traits. *Biometrika* **87**, 315–27.
- HUNG, M. & SWALLOW, W. H. (1999). Robustness of group testing in the estimation of proportions. *Biometrics* **55**, 231–7.
- HWANG, F. K. (1976). Group testing with a dilution effect. *Biometrika* **63**, 671–3.
- LE, C. T. (1981). A new estimator for infection rates using pools of variable size. *Am. J. Epidemiol.* **114**, 132–6.
- LITVAK, E., TU, X. M. & PAGANO, M. (1994). Screening for the presence of a disease by pooling sera samples. *J. Am. Statist. Assoc.* **89**, 424–34.
- RIDOUT, M. S. (1995). Three-stage designs for seed testing experiments. *J. R. Statist. Soc. C* **44**, 153–62.
- SOBEL, M. & ELASHOFF, R. (1975). Group testing with a new goal: estimation. *Biometrika* **62**, 181–93.
- SOBEL, M. & GROLL, P. A. (1959). Group testing to eliminate efficiently all defectives in a binomial sample. *AT&T Tech. J.* **38**, 1179–252.

- STEIN, C. (1945). A two-sample test for a linear hypothesis whose power is independent of the variance. *Ann. Math. Statist.* **16**, 243–58.
- TEBBS, J. M. & SWALLOW, W. H. (2003). Estimating ordered binomial proportions with the use of group testing. *Biometrika* **90**, 471–7.
- TU, X. M., LITVAK, E. & PAGANO, M. (1995). On the informativeness and accuracy of pooled testing in estimating prevalence of a rare disease: Application to HIV screening. *Biometrika* **82**, 287–9.
- UNITED STATES CENTERS FOR DISEASE CONTROL AND PREVENTION (2009). HIV prevalence estimates—United States, 2006. *J. Am. Med. Assoc.* **301**, 27–9.
- WEISS, S. H., GOEDERT, J. J., SARGADHARAN, M. G. & BODNER, A. J. (1985). The AIDS Seroepidemiology Collaborative working Group, Gallo RC, Blattner WA. Screening tests for HTLVI-III (AIDS Agent) antibodies: specificity, sensitivity, and applications. *J. Am. Med. Assoc.* **253**, 221–5.
- XIE, M., TATSUOKA, K., SACKS, J. & YOUNG, S. S. (2001). Group testing with blockers and synergism. *J. Am. Statist. Assoc.* **96**, 92–102.
- ZHANG, L., MUKHERJEE, B., GHOSH, M., GRUBER, S. & MORENO, V. (2008). Accounting for error due to misclassification of exposures in case–control studies of gene–environment interaction. *Statist. Med.* **27**, 2756–83.

[Received June 2010. Revised September 2011]