

## On Power and Sample Size Calculations for Likelihood Ratio Tests in Generalized Linear Models

Gwonen Shieh

Department of Management Science, National Chiao Tung University,  
Hsinchu, Taiwan 30050, Republic of China  
*email:gwshieh@cc.nctu.edu.tw*

**SUMMARY.** A direct extension of the approach described in Self, Mauritsen, and Ohara (1992, *Biometrics* 48, 31–39) for power and sample size calculations in generalized linear models is presented. The major feature of the proposed approach is that the modification accommodates both a finite and an infinite number of covariate configurations. Furthermore, for the approximation of the noncentrality of the noncentral chi-square distribution for the likelihood ratio statistic, a simplification is provided that not only reduces substantial computation but also maintains the accuracy. Simulation studies are conducted to assess the accuracy for various model configurations and covariate distributions.

**KEY WORDS:** Generalized linear models; Likelihood ratio test; Logistic regression; Noncentral chi-square; Poisson regression; Sample size; Score test; Statistical power.

### 1. Introduction

Generalized linear models were first introduced by Nelder and Wedderburn (1972) and are broadly applicable in almost all scientific fields. (See McCullagh and Nelder (1989) for further details.) The class of generalized linear models is specified by assuming independent scalar response variables  $Y_i$ ,  $i = 1, \dots, N$ , follow a probability distribution belonging to the exponential family of probability distributions with probability density of the form

$$\exp\{[Y\theta - b(\theta)]/a(\phi) + c(Y, \phi)\}. \quad (1.1)$$

The expected value  $E(Y) = \mu$  is related to the canonical parameter  $\theta$  by the function  $\mu = b'(\theta)$ , where  $b'$  denotes the first derivative of  $b$ . The link function  $g$  relates the linear predictors  $\eta$  to the mean response  $\eta = g(\mu)$ . The linear predictors can be written as

$$\eta = \mathbf{Z}^T \boldsymbol{\psi} + \mathbf{X}^T \boldsymbol{\lambda}, \quad (1.2)$$

where  $\mathbf{Z}(p \times 1)$  and  $\mathbf{X}(q \times 1)$  are vectors of covariates, and  $\boldsymbol{\psi}(p \times 1)$  and  $\boldsymbol{\lambda}(q \times 1)$  represent the corresponding unknown regression coefficients. The scale parameter  $\phi$  is assumed known. We wish to test the null hypothesis of  $H_0: \boldsymbol{\psi} = \boldsymbol{\psi}_0$  against the alternative hypothesis  $H_1: \boldsymbol{\psi} \neq \boldsymbol{\psi}_0$ , while treating  $\boldsymbol{\lambda}$  as a nuisance parameter.

For the purpose of power and sample size calculations within the framework of generalized linear models, two major tests have been proposed. They are the score test and likelihood ratio statistic developed by Self and Mauritsen (1988) and Self, Mauritsen, and Ohara (1992), respectively. However, these two approaches are limited to models where the number of covariate configurations is finite. This is somehow impracti-

cal since it is quite common for generalized linear models used in medical and clinical research to include continuous covariates as confounding factors, which have an infinite number of covariate configurations. For example, Whittemore (1981) illustrated the sample-size calculations for logistic regression with the problem of testing whether the incidence of coronary heart disease among white males aged 39–59 is related to their serum cholesterol and triglyceride levels. Also, previous studies indicate the joint distribution of cholesterol and log triglyceride is well presented by a bivariate normal distribution (see Hulley et al. (1980) for a thorough description of the analysis and other possible risk factors). In this case, to apply the approaches proposed by Self and Mauritsen (1988) and Self et al. (1992), one may apply a class grouping scheme over the range of covariate configurations. Such a strategy results in a categorical approximation of the true covariate distribution; hence, they are then still applicable with the consensus that the categorization is arbitrary. At first look, this seems to be a questionable approach because information about the actual serum cholesterol and triglyceride levels is thrown away. Furthermore, the interrelation between these two covariates may be distorted to some extent. Consequently, these two approaches do not fully exploit the distribution information about continuous covariates when it is available. More importantly, it is not clear how the results will be affected for utilizing a categorical approximation, not to mention that there is no unified rule for categorizing the covariates into a finite number of configurations.

In the present article, we generalize the Self et al. (1992) approach to accommodate covariates with an infinite number of configurations. With this natural modification, one can perform power and sample-size calculations in generalized linear

models with discrete and/or continuous covariates without any subjective or arbitrary class grouping process. In Section 2, the model and an approximation to the distribution of the likelihood ratio statistic are described. Section 3 provides the details of implementation. In Section 4, simulation studies are performed and results are presented that evaluate the adequacy of the proposed approach under various covariate distributions with an infinite number of configurations. Section 5 contains some remarks.

**2. Model and Approximation**

Consider a generalized linear model consisting of the response  $y_i$  and covariate  $(\mathbf{z}_i, \mathbf{x}_i)$  defined in (1.1) and (1.2), respectively, for  $i = 1, \dots, N$ . Assume  $(y_i, \mathbf{z}_i, \mathbf{x}_i)$  is a random sample from the joint distribution of  $(Y, \mathbf{Z}, \mathbf{X})$  with p.d.f.  $f(Y, \mathbf{Z}, \mathbf{X}) = f(Y | \mathbf{Z}, \mathbf{X}) \cdot f(\mathbf{Z}, \mathbf{X})$ , where  $f(Y | \mathbf{Z}, \mathbf{X})$  has the form defined in (1.1) and  $f(\mathbf{Z}, \mathbf{X})$  is the joint p.d.f. for  $\mathbf{Z}$  and  $\mathbf{X}$ . The form of  $f(\mathbf{Z}, \mathbf{X})$  is assumed to depend on none of the unknown parameters  $\psi$  and  $\lambda$ . The joint likelihood function of these  $N$  subjects is

$$L(\psi, \lambda) = \prod_{i=1}^N f(y_i, \mathbf{z}_i, \mathbf{x}_i) = \prod_{i=1}^N f(y_i | \mathbf{z}_i, \mathbf{x}_i) \cdot f(\mathbf{z}_i, \mathbf{x}_i).$$

It follows that the likelihood ratio statistic is  $2\{l(\hat{\psi}, \hat{\lambda}) - l(\psi_0, \hat{\lambda}_0)\}$ , where  $l(\psi, \lambda)$  is the log-likelihood function based on  $L(\psi, \lambda)$  and  $(\hat{\psi}, \hat{\lambda})$  and  $(\psi_0, \hat{\lambda}_0)$  are the maximum likelihood estimators of  $(\psi, \lambda)$  under the alternative and null models, respectively. The actual likelihood ratio test statistic is referred to its asymptotic distribution under the null hypothesis, which is a central chi-square distribution with  $p$  d.f. In order to perform power analysis, one needs to derive the distribution of the likelihood ratio statistic under the alternative hypothesis. Our formulation is analogous to that of Self et al. (1992). We approximate the distribution of the likelihood ratio statistic by a noncentral chi-square distribution with  $p$  d.f. The noncentrality parameter used in the approximation is computed by equating the expected value of a noncentral chi-square random variable to an approximation of the expected value of the likelihood ratio statistic. The expected value of the likelihood ratio statistic is approximated by the expected value of lead terms in an asymptotic expansion of the likelihood ratio statistic under the alternative hypothesis. As in Self et al. (1992), the likelihood ratio statistic is decomposed into three terms,

$$\begin{aligned} &2\{l(\hat{\psi}, \hat{\lambda}) - l(\psi_0, \hat{\lambda}_0)\} \\ &= 2\{l(\hat{\psi}, \hat{\lambda}) - l(\psi, \lambda)\} - 2\{l(\psi_0, \hat{\lambda}_0) - l(\psi_0, \lambda_0^*)\} \\ &\quad + 2\{l(\psi, \lambda) - l(\psi_0, \lambda_0^*)\}, \end{aligned} \tag{2.1}$$

where  $\lambda_0^*$  is the limiting value of  $\hat{\lambda}_0$  as described in Self and Mauritsen (1988). To approximate the expected value of the first term in (2.1), only the lead term in Cordeiro's (1983) expansion for generalized linear models is retained, and it results in a value of  $p + q$ .

The approximation of the second term is more troublesome because none of the expansions in Bartlett (1953), Lawley (1956), and Cordeiro (1983) are performed about the parameter  $(\psi_0, \lambda_0^*)$ . However, the expected value of the first term in the expansion is the trace of two  $q \times q$  matrices,  $\Sigma^{-1}$  and

$\Xi$ ,  $\text{tr}(\Sigma^{-1}\Xi)$ , where

$$\begin{aligned} \Sigma &= E \left[ -\frac{\partial^2 l(\psi, \lambda)}{\partial \lambda^2} \Big|_{(\psi_0, \lambda_0^*)} \right] \\ &= N \cdot E_{(\mathbf{Z}, \mathbf{X})} \left[ a^{-1}(\phi) \left\{ b''(\theta^*) \left( \frac{\partial \theta^*}{\partial \eta^*} \right)^2 \right. \right. \\ &\quad \left. \left. - [b'(\theta) - b'(\theta^*)] \left( \frac{\partial^2 \theta^*}{\partial \eta^{*2}} \right) \right\} \right] \end{aligned}$$

and

$$\begin{aligned} \Xi &= E \left[ \left\{ \frac{\partial l(\psi, \lambda)}{\partial \lambda} \Big|_{(\psi_0, \lambda_0^*)} \right\} \left\{ \frac{\partial l(\psi, \lambda)}{\partial \lambda} \Big|_{(\psi_0, \lambda_0^*)} \right\}^T \right] \\ &= N \cdot E_{(\mathbf{Z}, \mathbf{X})} \left[ a^{-1}(\phi) b''(\theta) \left( \frac{\partial \theta^*}{\partial \eta^*} \right)^2 \mathbf{X} \mathbf{X}^T \right], \end{aligned}$$

$E[\cdot]$  denotes the expectation taken with respect to the joint distribution of  $(Y_1, \dots, Y_N, \mathbf{Z}, \mathbf{X})$  under the alternative hypothesis with true parameter values  $(\psi, \lambda)$ , and  $E_{(\mathbf{Z}, \mathbf{X})}[\cdot]$  denotes the expectation taken with respect to the joint distribution of  $(\mathbf{Z}, \mathbf{X})$ , which is not dependent on the value of  $(\psi, \lambda)$ . Also,  $b''$  denotes the second derivative of  $b$ ,  $\theta$  and  $\theta^*$  denote the canonical parameter values evaluated at  $(\psi, \lambda)$  and  $(\psi_0, \lambda_0^*)$ , respectively, and  $\eta^*$  denotes the linear predictor evaluated at  $(\psi_0, \lambda_0^*)$ . It was found in Self et al. (1992) and Shieh and O'Brien (1998) that the value of  $\text{tr}(\Sigma^{-1}\Xi)$  is very close to  $q$  for certain parameter values and discrete covariate distributions in several generalized linear models. This phenomenon is fortified here from the numerical results in this paper. We found that there is essentially no practical difference in the adequacy for power and sample size calculations by replacing it with  $q$ .

The third term does not involve any maximum likelihood estimators of  $(\psi, \lambda)$ . Its expectation can be written as  $N\Delta^*$ , where

$$\Delta^* = E_{(\mathbf{Z}, \mathbf{X})} \left[ 2a^{-1}(\phi) \{ b'(\theta)[\theta - \theta^*] - [b(\theta) - b(\theta^*)] \} \right]. \tag{2.2}$$

By equating the expected value of a noncentral chi-square random variable with  $p$  d.f. and noncentrality  $\gamma$ , namely  $p + \gamma$ , to the respective expected value approximations of (2.1) derived above, which is  $(p + q) - q + N\Delta^*$ . This leads to an estimate of noncentrality, denoted by  $\gamma_N$ , of the proposed noncentral chi-square distribution for the likelihood ratio statistic under the alternative hypothesis, i.e.,  $\gamma_N = N\Delta^*$ . The subscript  $N$  of  $\gamma_N$  represents its dependency on sample size  $N$ . Hence, for given parameter values  $(\psi, \lambda)$  and covariate distribution  $f(\mathbf{Z}, \mathbf{X})$ , the actual likelihood ratio statistic of sample size  $N$  under the alternative hypothesis is performed by referring it to a noncentral chi-square distribution with  $p$  d.f. and noncentrality  $N\Delta^*$ . Our approach differs from Self et al. (1992) in two respects. First, they proposed considering only categorical covariates, which are restricted to have a finite number of configurations such as Bernoulli or multinomial distributions. This is naturally extended here since the joint distribution of covariates  $(\mathbf{Z}, \mathbf{X})$  could be either discrete or continuous with an infinite number of configurations, e.g., Poisson and normal distributions. The expression of  $\Delta^*$  in (2.2) subsumes Self et al.'s equation (2.3) as a special case when

the joint distribution of  $(\mathbf{Z}, \mathbf{X})$  is categorical with probabilities  $\pi_j$ ,  $j = 1, \dots, m$ . Second, their noncentrality estimate is  $N\Delta^* + q - \text{tr}(\Sigma^{-1}\Xi)$ , whereas ours is simply  $N\Delta^*$ .

**3. Implementation**

In this section, we shall describe the necessary steps to implement the proposed approach. For a generalized linear model with specified parameter values  $(\psi, \lambda)$  and chosen covariate distribution  $f(\mathbf{Z}, \mathbf{X})$ , the sample size needed to test hypothesis  $H_0: \psi = \psi_0$  with specified significance level  $\alpha$  and power  $1 - \beta$  against the alternative  $H_1: \psi \neq \psi_0$  is computed as follows. First, find the  $100(1 - \alpha)$ th percentile of a central chi-square distribution with  $p$  d.f., denoted by  $\chi^2_{p,1-\alpha}$ . Next, find the noncentrality  $\gamma_N$  of a noncentral chi-square distribution with  $p$  d.f. such that the  $100 \cdot \beta$ th percentile, denoted by  $\chi^2_{p,\beta}(\gamma_N)$ , equals  $\chi^2_{p,1-\alpha}$ . Then the sample size estimate  $N$  is computed as  $\gamma_N/\Delta^*$ , where  $\Delta^*$  is as defined in (2.2). For continuous covariate distributions, numerical integration is needed to carry out the expectation for  $\Delta^*$ .

**4. Simulation Studies**

Simulation studies are performed for evaluating the accuracy of the proposed approach for logistic regression models and Poisson regression models with an infinite number of covariate configurations. For illustrative purposes, we will restrict our attention to the logistic regression models here.

Two sets of linear predictors of the form  $\eta = \lambda_I + Z\psi$  and  $\eta = \lambda_I + X\lambda_S + Z\psi$  are examined. In the case of the simple linear predictor,  $\eta = \lambda_I + Z\psi$ , we consider normal, double exponential, exponential, and Poisson distributions for covariate  $Z$ . For the second linear predictor,  $\eta = \lambda_I + X\lambda_S + Z\psi$ , the joint distribution of  $(Z, X)$  is of the form  $f(Z, X) = f(X | Z)f(Z)$ . We assume  $Z$  is a Bernoulli covariate with  $p(Z = 1) = \pi$  for  $\pi = 0.1, 0.5, \text{ and } 0.9$ . The conditional distribution of  $X$  given  $Z$ , denoted by  $[X | Z]$ , is  $[X | Z = 1] \equiv [X | Z = 0] + d$ , where  $[X | Z = 0]$  is a standardized version of a normal, double exponential, exponential, or Poisson with a mean of 10 random variables and  $d$  is described in the footnote of Table 2.

The parameter of interest,  $\psi$ , is taken to be  $\log(2)$  and

$\log(5)$  for the two linear predictors, respectively. The confounding parameter  $\lambda_S$  in the second model is set as  $\log(2)$ . The intercept parameter,  $\lambda_I$ , is chosen to satisfy different values of overall response probability  $\bar{\mu} = E_{(\mathbf{Z}, \mathbf{X})}[\exp(\eta)/\{1 + \exp(\eta)\}]$ .

For a given model, covariate distribution, regression coefficients, and overall response probability, the estimates of sample size required for testing  $H_0: \psi = 0$  against the alternative hypothesis  $H_1: \psi \neq 0$  with significance level 0.05 and power 0.80, 0.90, and 0.95 are calculated. Due to the rounding of sample sizes, the precise nominal powers are not exactly 0.80, 0.90, and 0.95. They are recalculated with the inversion of the proposed formula discussed in Section 3.

Estimates of actual power associated with a given sample size and model configurations are then computed through Monte Carlo simulation based on 5000 replicate data sets. The adequacy of the sample size formula is determined by the difference between the estimated power and nominal power specified above. All calculations are performed using programs written with SAS/IML (SAS, 1989).

The results of the simulation studies are presented in Tables 1 and 2. Table 1 contains results for the simple linear predictor, while Table 2 contains results of the multiple linear predictor for  $p(Z = 1) = 0.5$ . In general, the sample size needed to achieve the significance level and power is larger for overall response probability  $\bar{\mu} = 0.02$  than for  $\bar{\mu} = 0.15$ . However, for most occasions, the absolute errors of overall response probability  $\bar{\mu} = 0.02$  are larger than those of  $\bar{\mu} = 0.15$ . Hence, the proposed method is comparatively more accurate for larger overall response probability. Generally, it maintains the accuracy within a reasonable range of nominal power for both cases. The simulation results for logistic regression in different settings and for Poisson regression models also suggest the proposed method performs well and may be of practical use; however, they are not reported here.

**5. Discussion**

We propose in this article an approach for sample size and power calculations in generalized linear models. This approach is a direct extension of the work by Self et al. (1992) to ac-

**Table 1**

*Calculated sample sizes and estimates of actual power for the logistic regression model with linear predictor  $\eta^a = \lambda_I + Z\psi$*

	Distribution of $Z^b$											
	Normal		Double exponential		Exponential		Poisson					
$\bar{\mu} = 0.02$												
Sample size	849	1136	1405	668	894	1105	368	492	608	567	758	938
Nominal power	.8003	.9001	.9501	.8004	.9002	.9500	.8009	.9003	.9500	.8006	.9001	.9502
Estimated power	.7774	.8818	.9406	.7504	.8602	.9168	.7190	.8328	.8888	.7586	.8534	.9144
Error	-.0229	-.0183	-.0095	-.0500	-.0400	-.0332	-.0819	-.0675	-.0612	-.0420	-.0467	-.0358
$\bar{\mu} = 0.15$												
Sample size	141	189	233	139	186	230	101	136	168	113	152	187
Nominal power	.8011	.9012	.9502	.8018	.9012	.9507	.8002	.9018	.9306	.8002	.9015	.9500
Estimated power	.7930	.9016	.9470	.7968	.8908	.9444	.7742	.8768	.9306	.7908	.8890	.9338
Error	-.0081	.0004	-.0032	-.0050	-.0104	-.0063	-.0260	-.0250	-.0203	-.0094	-.0125	-.0162

<sup>a</sup>  $\psi = \log(2)$  and  $E_{(Z)}[\exp(\eta)/\{1 + \exp(\eta)\}] = \bar{\mu}$ .

<sup>b</sup> The distribution of  $Z$  is standardized to have mean zero and variance one.

**Table 2**  
 Calculated sample sizes and estimates of actual power for the logistic regression model with linear predictor  $\eta^a = \lambda_I + X\lambda_S + Z\psi$ , where  $Z$  has a Bernoulli distribution ( $\pi = 0.5$ )

	Distribution of $[X   Z]^b$											
	Normal		Double exponential		Exponential		Poisson					
$\bar{\mu} = 0.02$												
Sample size	2250	3011	3724	1552	2078	2569	1460	1955	2417	648	867	1073
Nominal power	.8002	.9000	.9500	.8001	.9001	.9500	.8001	.9001	.9500	.8003	.9000	.9502
Estimated power	.8354	.9280	.9680	.8420	.9252	.9744	.8326	.9328	.9728	.8372	.9218	.9692
Error	.0352	.0280	.0180	.0419	.0251	.0244	.0325	.0327	.0228	.0369	.0218	.0190
$\bar{\mu} = 0.15$												
Sample size	272	364	450	207	276	342	194	260	322	129	172	213
Nominal power	.8004	.9002	.9501	.8017	.9001	.9504	.8000	.9003	.9505	.8027	.9008	.9508
Estimated power	.8262	.9150	.9610	.8168	.9114	.9582	.8226	.9082	.9576	.8092	.9072	.9540
Error	.0258	.0148	.0109	.0151	.0113	.0078	.0226	.0079	.0071	.0065	.0064	.0032

<sup>a</sup>  $\psi = \log(5)$ ,  $\lambda_S = \log(2)$ , and  $E_{(Z,X)}[\exp(\eta)/\{1 + \exp(\eta)\}] = \bar{\mu}$ .

<sup>b</sup> The distribution  $[X | Z = 0]$  is standardized to have mean zero and variance one. The distribution of  $[X | Z = 1] \equiv [X | Z = 0] + d$ , where  $d$  is 1.6832, 1.2958, 1.3863, and  $5/(10)^{1/2}$  for normal, double exponential, exponential, and Poisson distributions, respectively.

commodate covariate distributions with an infinite number of configurations. Their approach is restricted to the generalized linear models with a finite number of covariate configurations such as Bernoulli and multinomial distributions. Furthermore, we modify the approximation of the noncentrality parameter in a noncentral chi-square distribution of the likelihood ratio statistic. This simple structure permits computational simplifications and maintains great accuracy based on the simulation results for different settings of logistic regression and Poisson regression models.

For generalized linear models with continuous covariates of natural interval and ratio measurement scales, some researchers may prefer to work with a categorical approximation by grouping the range of covariate values into finite intervals and then choosing representative class values (usually the class midpoints) and proportions for each class. This process will make the Self et al. (1992) approach still applicable for models with an infinite number of covariate configurations. However, there is no consensus in determining the covariate distribution approximation in terms of numbers of classes, the choices of class boundaries, and the class representative values. Consequently, one classification scheme may perform well for some cases but do poorly for others. Along with the proposed approach, we have simultaneously evaluated two different classification schemes for each of the four covariate distributions in the simulation studies. The results indicate that the proposed approach outperforms those two with categorical approximations of covariate distributions for most of the cases that we have considered. Therefore, when the covariate distributions are available, one should incorporate such information into the sample size calculations instead of their categorical approximations. However, as pointed out by the referee, the latter may be more robust when the distribution information about covariates is not accurately known. In such cases, one may try several different settings of finite configurations to provide guidance about the sample sizes required for the study.

ACKNOWLEDGEMENTS

The author wishes to thank the associate editor and referees for their helpful comments and Dr Ralph O'Brien for enlightening discussions.

RÉSUMÉ

Cet article présente une extension directe de l'approche décrite dans Self, Mauritsen and Ohara (1992, *Biometrics* **48**, 31-39) pour les calculs de puissance et de taille d'échantillon pour les modèles linéaires généralisés. L'apport majeur de l'approche proposée est la modification qui permet aussi bien un nombre fini qu'infini de configurations de covariables. De plus il est aussi proposé une simplification pour l'approximation du paramètre de non centralité de la distribution du chi-deux non centré, approximation qui non seulement réduit les calculs de manière appréciable, mais aussi conserve la précision. Des études de simulation sont faites pour évaluer cette précision pour différentes configurations de modèles et de distributions des covariables.

REFERENCES

Bartlett, M. A. (1953). Approximate confidence intervals. II. More than one unknown parameter. *Biometrika* **40**, 306-317.

Cordeiro, G. M. (1983). Improved likelihood ratio statistics for generalized linear models. *Journal of the Royal Statistical Society, Series B* **45**, 404-413.

Hulley, S. B., Rosenman, R. A., Bawol, R., and Brand, R. J. (1980). Epidemiology as a guide to clinical decisions: The association between triglyceride and coronary heart disease. *New England Journal of Medicine* **302**, 1383-1389.

Lawley, D. N. (1956). A general method for approximating the distribution of likelihood ratio criteria. *Biometrika* **43**, 295-303.

McCullagh, P. and Nelder, J. A. (1989). *Generalized Linear Models*, 2nd edition. London: Chapman and Hall.

- Nelder, J. A. and Wedderburn, R. W. M. (1972). Generalized linear models. *Journal of the Royal Statistical Society, Series A* **135**, 370–384.
- SAS. (1989). *SAS/IML Software: Usage and Reference*, Version 6. Cary, North Carolina: SAS Institute.
- Self, S. G. and Mauritsen, R. H. (1988). Power/sample size calculations for generalized linear models. *Biometrics* **44**, 79–86.
- Self, S. G., Mauritsen, R. H., and Ohara, J. (1992). Power calculations for likelihood ratio tests in generalized linear models. *Biometrics* **48**, 31–39.
- Shieh, G. and O'Brien, R. G. (1998). A simpler method to compute power for likelihood ratio tests in generalized linear models. Paper presented at the Annual Joint Statistical Meetings of the American Statistical Association, Dallas, Texas.
- Whittemore, A. S. (1981). Sample size for logistic regression with small response probability. *Journal of the American Statistical Association* **76**, 27–32.

*Received August 1999. Revised January 2000.*

*Accepted March 2000.*