

Tilburg University

Case-Control Studies with Contaminated Controls

Imbens, G.W.; Lancaster, T.

Publication date:
1993

[Link to publication in Tilburg University Research Portal](#)

Citation for published version (APA):
Imbens, G. W., & Lancaster, T. (1993). *Case-Control Studies with Contaminated Controls*. (CentER Discussion Paper; Vol. 1993-7). CentER.

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

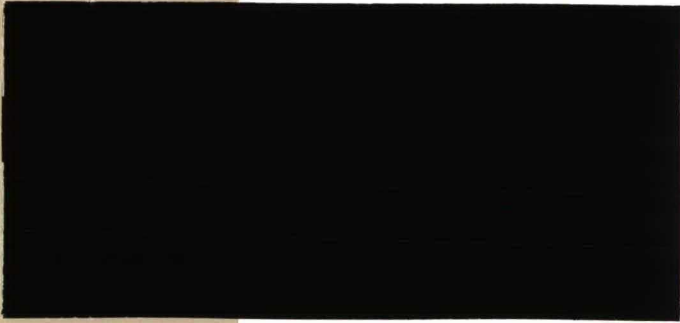
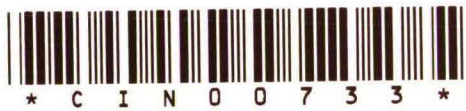
Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

ECO
CBM
R R
8414
1993
7

Center
for
Economic Research

Discussion paper



No 9307

× CASE-CONTROL STUDIES
WITH CONTAMINATED CONTROLS

by Guido Imbens and
Tony Lancaster

January 1993

I 0924-7815

CASE-CONTROL STUDIES
WITH CONTAMINATED CONTROLS†

GUIDO IMBENS – HARVARD UNIVERSITY*
TONY LANCASTER – BROWN UNIVERSITY** AND HARVARD UNIVERSITY

JULY 1992

†We are grateful for financial support from the National Science Foundation under grant SES 9122477 and to CentER, Tilburg University for hospitality.

* mailing address: Department of Economics, Harvard University, Cambridge, MA 02138

** mailing address: Department of Economics, Brown University, Providence, RI 02912

1. INTRODUCTION.

There is a significant body of literature in statistics and econometrics dealing with discrete response models under various types of non-random sampling. Such sampling schemes might reduce the cost of the study, particularly if one of the responses is rare. A leading case is case-control, retrospective, choice-based or response-based sampling. In the simplest example the researcher has two samples, one containing observations with response $y = 1$ (the cases), and the second containing observations with response $y = 0$ (the controls). In both samples we observe the attributes x for all observations. When the model for the conditional probabilities of the choices given the covariates is of logit form it has long been known that the investigator can proceed as though the data were obtained by random sampling so far as estimation of the covariate coefficients is concerned; see for example Prentice and Pyke(1979). For the general case Manski and Lerman (1977) proposed a weighted maximum likelihood estimator. Cosslett (1981) and Imbens (1990) proposed efficient solutions to the general estimation problem.

A case that has not received as much attention, and one that is not covered by the general sampling schemes in Hsieh, Manski and McFadden (1985) and Imbens (1990) is that where the second sample is a random sample *from the whole population* with only the attributes or covariate values and not the responses, observed. The second sample, that formed the control group in case-control sampling, now consists of an unknown mixture of cases and controls. Such a situation might occur if the researcher obtains a sample of observations with a particular response, for example being a labor force participant or being unemployed, and wishes, possibly for reasons of economy, to compare them with a random sample from a very different source in which the particular response was not measured. We describe this set up as one of contaminated controls, following the usage of Heckman and Robb(1984). Neither sample in itself identifies the parameters of the conditional response probability but the combination of cases and contaminated controls might do so.

This paper deals with efficient estimation of parametric discrete choice using samples of this type. In section 2 we discuss identifiability of choice models under contaminated sampling and point out that the choice model is nonparametrically identified if the marginal

probabilities of the choices are known to the investigator. In section 3 we give an efficient generalized method of moments (GMM) estimator for the case in which the marginal probabilities are unknown. The estimator is identical to a constrained maximum likelihood estimator when the covariates have a multinomial distribution with known support. In section 4 we give an efficient GMM estimator for the case in which the marginal probabilities are known. This estimator is asymptotically equivalent to a constrained maximum likelihood estimator when the covariates are multinomial. The estimator proposed in section 3 achieves the semiparametric efficiency bound as defined by Chamberlain (1987) or Begun et al (1983). The problem is semiparametric because of the appearance in the likelihood of the unknown population covariate distribution.

In section 5 we discuss the case in which the choice model is logit and the marginal probabilities are known. This case has been considered by Steinberg and Cardell (1991) who have given a consistent estimator of the logit parameters. Section 6 reports a small Monte Carlo study of the estimators.

2. THE MODEL AND ITS IDENTIFIABILITY

Let y be a binary random response variable, equal to 0 or 1, and x a vector of attributes. In the population the distribution function of x is $F(x)$ which is unknown. We will assume that the conditional probability of $y = 1$ given x in the population is equal to $pr(y = 1|x) = P(x; \beta)$ where $P(\cdot; \cdot)$ is a known function and β an unknown parameter. Finally, we define q to be the marginal probability of choice 1 in the population, $q = \int P(x; \beta) dF(x)$.

The sampling scheme is that two independent random samples of sizes N_1 and N_0 are available. The first is drawn from the subset of the population who made choice 1 and the covariate is observed; the second is drawn from the whole population with only the covariate observed. We let s denote a binary stratum indicator, taking the value 1 if an observation is drawn from the sub-population who made choice 1, and 0 if it was drawn from the whole population.

An observation from stratum 1 has probability $p(x|y = 1) = P(x)f(x)/q$; an observation from stratum 0 has probability $f(x)$. If we knew these probabilities we could determine the function $P(x)/q$ for all values of x with positive probability. This function is therefore non-parametrically identified. It follows that the relative probabilities $P(x)/P(x_0)$ are identified. This contrasts with standard case-control sampling which identifies the relative odds, $P(x)/(1 - P(x)) \div P(x_0)/(1 - P(x_0))$.

If q is also known then clearly $P(x)$ is identifiable. Alternatively, if the parametric form of $P(x; \beta)$ is known then β can generally be deduced from knowledge of the function $P(x)/q$ for a sufficiently large set of values of x . In this case $P(x)$ is parametrically identifiable. In this paper we shall consider parametric models for $P(x)$ with and without prior knowledge of q . When q is known $P(x)$ is parametrically overidentified.

3. EFFICIENT ESTIMATION.

In this section we will propose an estimator for the parameters of the conditional choice probability function $P(x; \beta)$. This function $P(x)$ will be assumed known up to a finite parameter vector β and there is no prior knowledge of the marginal probability q . In section 4 we shall show how to take account of prior information such as knowledge of q .

To derive this estimator we will assume initially that the regressors x have a discrete distribution with unknown probabilities λ_l on $L + 1$ known points of support, x^l . This allows us to use standard maximum likelihood theory, and to derive an efficient estimator for that case. This estimator does not depend on either the number or the location of points of support of the covariate distribution that do not appear in the sample. We then show that this estimator is asymptotically semiparametrically efficient.

It is convenient, first of all, to enlarge the model. We do this by supposing that the sample sizes were determined by a sequence of Bernoulli trials with unknown parameter h . Thus the data is provided by repeatedly conducting such trials; if a success occurs we randomly sample from the subpopulation who made choice 1; if a failure, we randomly sample from the whole population. This procedure is repeated N times. The population is assumed sufficiently large that the probability of overlap between the sampled individuals is zero. A consequence of this enlargement is that the sample now constitutes N independently and identically distribution realisations from the joint distribution of stratum and covariate $g(s, x) = (hPf/q)^s((1-h)f)^{1-s}$. The quantity h will be treated as an unknown parameter. Its maximum likelihood estimator will be the sample fraction of observations from stratum 1, N_1/N . As long as h is functionally independent of β , N_1/N is ancillary and the asymptotic distribution of the ML estimator of β is independent of that of h .

If $N = N_1 + N_0$ is the total number of observations the log likelihood is

$$L(\beta, h, \lambda) = \sum_{n=1}^N [s_n \log[P_n(\beta)f_n(\lambda)/q(\beta, \lambda)] + (1 - s_n) \log f_n(\lambda)] \\ + N_1 \log h + N_0 \log(1 - h), \quad (3.1)$$

where $f_n(\lambda) = f(x_n; \lambda)$ and $P_n(\beta) = P(x_n; \beta)$. Since L involves β, λ in a rather awkward way because of the term in q it is convenient to reparametrize. The following transformation

changes the log likelihood into the form that would arise under a random sampling scheme in which there exists a conditional distribution and a marginal distribution each depending on distinct sets of parameters.

Define

$$R_1(x; \beta, q, h) = \frac{(h/q)P(x; \beta)}{(h/q)P(x; \beta) + 1 - h}, \quad R_0 = 1 - R_1, \quad (3.2)$$

$$g(x) = [(h/q)P(x; \beta) + 1 - h]f(x).$$

R_1 is the conditional probability that an observation comes from stratum 1 given the covariate and the sampling scheme. The distribution $g(x)$, which is also multinomial with parameters $\pi_l = [(h/q)P(x^l; \beta) + 1 - h]\lambda_l$ on the same points of support as $f(x)$, is the covariate distribution induced by the sampling scheme. Then L may be rewritten as

$$L(\beta, q, h, \pi) = \sum_{n=1}^N [s_n \log R_{1n}(\beta, q, h) + (1 - s_n) \log R_{0n}(\beta, q, h)]$$

$$+ \sum_{n=1}^N \log g_n(\pi)$$

$$= L_1(\beta, q, h) + L_2(\pi) \quad (3.3)$$

We can regard L as a function of the parameters β, q, h, π , where these parameters are subject to the constraint that $q = \int P(x; \beta) dF(x; \lambda)$ which may be rewritten in terms of the new parametrization as

$$h = \int R_1(x; \beta, q, h) dG(x; \pi). \quad (3.4)$$

We now give the ML estimator of β, q, h, π . Let a hat denote an estimator which maximizes L without imposing the restriction (3.4). Then $\hat{\pi}_l = n_l/N$ for all l where n_l is the sample number of observations which have covariate value x^l . At this solution for π the constraint, (3.4), becomes

$$h = N^{-1} \sum_{n=1}^N R_{1n}(\beta, q, h). \quad (3.5)$$

Next consider the β , q and h likelihood equations from L_1 .

$$\frac{\partial L_1}{\partial \beta} = \sum_{n=1}^N p'_{\beta n}(s_n - R_{1n}(\beta, q, h))/P_n = 0 \quad (3.6)$$

$$\frac{\partial L_1}{\partial q} = -(1/q) \sum_{n=1}^N (s_n - R_{1n}(\beta, q, h)) = 0 \quad (3.7)$$

$$\frac{\partial L_1}{\partial h} = (1/\bar{h}) \sum_{n=1}^N (s_n - R_{1n}(\beta, q, h)) = 0. \quad (3.8)$$

Here $p_{\beta n} = \partial P_n / \partial \beta$ of order $1 \times K$ where K is the dimension of β and $\bar{h} = h(1 - h)$.

Let $\hat{\beta}, \hat{q}$ solve (3.6) and (3.7) with $h = \hat{h} = N_1/N$. Then $\hat{\beta}, \hat{q}, \hat{h}$ solve (3.6), (3.7), (3.8) and they also satisfy the constraint (3.5) which may be written $N^{-1} \sum (s_n - R_{1n}(\hat{\beta}, \hat{q}, \hat{h})) = 0$. Hence the constrained ML estimator of β, q can be found by maximising $L_1(\beta, q, \hat{h})$ with respect to variation in β, q . Since L_1 is just a random sampling binary choice log likelihood this is an essentially simple computation.

The above derivation gives $\hat{\beta}$ as a constrained ML estimator after a parameter transformation. It may also be given a generalized method of moments (GMM) interpretation.¹ Consider the generalized moments

$$\begin{aligned} \psi_1(\beta, q, h, s, x) &= p'_\beta(x; \beta)(s - R_1(x; \beta, q, h))/P(x; \beta) \\ \psi_2(\beta, q, h, s, x) &= -(1/q)(s - R_1(x; \beta, q, h)) \\ \psi_3(\beta, h, q, s, x) &= q - P(x; \beta)/[(h/q)P(x; \beta) + 1 - h] \propto h - R_1(x; \beta, q, h). \end{aligned} \quad (3.9)$$

The moments ψ_1, ψ_2 are the single observation scores for β, q from the log likelihood L_1 , (3.3). In the form $q - P/[(h/q)P + 1 - h]$ the moment ψ_3 is just the definitional relation between marginal, q , and conditional, $P(x)$, choice probabilities after allowing for the fact that the covariate distribution induced by the sampling scheme is not $f(x)$, but $g(x) = f(x)/[(h/q)P + 1 - h]$. In the form $h - R_1$ the moment ψ_3 is the single observation version of the constraint (3.4). These moments have mean zero at the true parameter point.

¹See Hansen(1982), Manski(1988).

Equating their sample analogues to zero gives $\hat{\beta}, \hat{q}, \hat{h}$ which are then GMM estimates. Thus the asymptotic distribution of the estimator may be found equivalently from GMM theory or from constrained ML theory. The former is rather simpler since we do not have to consider the estimation of π . Moreover note that these are valid moments whether the distribution of x is discrete or continuous so they do not hinge on the assumption of a discrete covariate with known support.

Theorem 1. Let $\delta = (\beta, q, h)$ and $\psi = (\psi_1, \psi_2, \psi_3)$ where $\psi_3 = h - R_1(x; \beta, q, h)$. Under regularity conditions, the solution, $\hat{\delta}$ to $\sum_{n=1}^N \psi_n(\hat{\delta}) = 0$ is a consistent estimator for δ^* and $\sqrt{N}(\hat{\delta} - \delta^*) \rightarrow \mathcal{N}(0, V)$ where

$$V = \Gamma^{-1} \Delta (\Gamma')^{-1}, \quad \Delta = \mathcal{E}[\psi(\delta) \cdot \psi(\delta)']_{\delta=\delta^*}, \quad \Gamma = \mathcal{E} \left[\frac{\partial \psi}{\partial \delta \delta'} \right]_{\delta=\delta^*}.$$

An asterisk denotes the true value. The above covariance matrix is the semiparametric efficiency bound of Chamberlain(1987) or Begun, Hall, Huang and Wellner(1984). Proof: see appendix.

An explicit form for the asymptotic covariance matrix of $\hat{\beta}, \hat{q}$ is as follows. Let

$$\begin{aligned} \Delta_{11} &= \mathcal{E}(p'_\beta \bar{R} p_\beta / P^2); & \Delta_{12} &= -(1/q) \mathcal{E}(p'_\beta \bar{R} / P) \\ \Delta_{22} &= (1/q^2) \mathcal{E}(\bar{R}); & \Delta_{33} &= \bar{h} - \mathcal{E}(\bar{R}) \end{aligned} \quad (3.10)$$

which are the non-zero elements of Δ . Here $\bar{R} = R_1(1 - R_1)$ and the expectation is with respect to $g(x)$, defined in (3.2). Then the limiting covariance matrix of $\hat{\beta}, \hat{q}$ is

$$V(\hat{\beta}, \hat{q}) = \Delta_1^{-1} - \begin{pmatrix} 0 & 0 \\ 0 & q^2/\bar{h} \end{pmatrix} \text{ where } \Delta_1 = \begin{pmatrix} \Delta_{11} & \Delta_{12} \\ \Delta_{21} & \Delta_{22} \end{pmatrix}. \quad (3.11)$$

The variance of \hat{h} is \bar{h} and it is distributed independently of $\hat{\beta}, \hat{q}$.

We see that the covariance matrix of $\hat{\beta}$ can be found from the upper left submatrix of Δ_1^{-1} which is the inverse information matrix for β, q from L_1 . This means that (a) an efficient estimate of β can be found by maximizing the binary choice log likelihood, L_1 , with respect to β, q with h replaced by N_1/N , and (b) the standard inverse information matrix estimate of the $\hat{\beta}, \hat{q}$ covariance matrix will give the correct standard errors for $\hat{\beta}$ (though not for \hat{q} .)

4. EFFICIENT ESTIMATION WITH KNOWN q

Suppose that extra sample information provides the numerical value of the marginal choice probability, q^* . One way of proceeding is to maximize the log likelihood (3.3) subject to the constraint provided by knowledge of q^* . The log likelihood becomes

$$\begin{aligned} L(\beta, h, \pi) &= \sum_{n=1}^N [s_n \log R_{1n}(\beta, q^*, h) + (1 - s_n) \log R_{0n}(\beta, q^*, h)] \\ &\quad + \sum_{n=1}^N \log g_n(\pi) \\ &= L_1(\beta, h) + L_2(\pi) \end{aligned} \quad (4.1)$$

The constraint relating β, h, π is $q^* = \int P(x; \beta) dF(x; \lambda)$ which is equivalent to

$$h = \int R_1(x; \beta, h) dG(x; \pi). \quad (4.2)$$

Here, $R_1 = (h/q^*)P/[(h/q^*)P + 1 - h]$. The ML estimator of β, h, π maximizes (4.1) subject to (4.2). Unlike the case in which q was unknown it is no longer true that the unconstrained ML estimator satisfies the constraint, so this simplification no longer applies. But a constrained optimization can be avoided if we adopt a Generalized Method of Moments approach.

Consider the moments ψ with q replaced q^* . These are

$$\begin{aligned} \psi_1(\beta, q^*, h, s, x) &= p'_\beta(x; \beta)(s - R_1(x; \beta, q^*, h))/P(x; \beta) \\ \psi_2(\beta, q^*, h, s, x) &= -(1/q^*)(s - R_1(x; \beta, q^*, h)) \\ \psi_3(\beta, q^*, h, s, x) &= h - R_1(x; \beta, q^*, h). \end{aligned} \quad (4.3)$$

The covariance matrix of these moments is Δ whose elements were given in (3.10). Then

Theorem 2.

Let $\psi_n = \psi(\beta, q^*, h, s_n, x_n)$; $\delta = (\beta, h)$, $\Delta = \mathcal{E}(\psi\psi')$ and $\Gamma = \mathcal{E}(\partial\psi/\partial\delta)$. Γ is now a submatrix of the Γ of theorem 1 — the column corresponding to q has been deleted. Finally, let $\hat{\delta}$ minimize

$$\sum_{n=1}^N \psi_n(\delta)\Delta^{-1}\psi_n(\delta).$$

Then $\sqrt{N}(\hat{\delta} - \delta^*) \rightarrow \mathcal{N}(0, V)$ where

$$V = [\Gamma' \Delta^{-1} \Gamma]^{-1}.$$

This covariance matrix is the same as that of the estimator of β, h which maximizes (4.1) subject to (4.2), under the usual regularity conditions. Thus the GMM estimator is asymptotically equivalent to the ML estimator and is efficient when the covariate is discrete with known points of support. We conjecture that $\hat{\beta}$ is also semiparametrically efficient.

Notice the simplicity of the GMM procedure. It avoids estimation of the covariate distribution; it avoids a constrained optimization problem; and it is a procedure that can be applied without any restrictive assumption about the covariate distribution.

An explicit form for the asymptotic covariance matrix of $\hat{\beta}$ is²

$$V(\hat{\beta}) = \Delta_{11}^{-1} - \Delta_{11}^{-1} \Delta_{12} [\Delta_{21} \Delta_{11}^{-1} \Delta_{21} + (\bar{h}/q^2) - \Delta_{22}]^{-1} \Delta_{21} \Delta_{11}^{-1} \quad (4.4)$$

The corresponding expression when q is not known is found from (3.11) to be

$$V(\hat{\beta}) = \Delta_{11}^{-1} - \Delta_{11}^{-1} \Delta_{12} [\Delta_{21} \Delta_{11}^{-1} \Delta_{21} - \Delta_{22}]^{-1} \Delta_{21} \Delta_{11}^{-1} \quad (4.5)$$

The feasible form of the estimator will require an initial consistent estimate of δ in order to estimate the covariance matrix Δ . This might be provided by the estimator which solves

$$\sum_{n=1}^N \psi_1(\beta, q^*, \hat{h}, s_n, x_n) = 0 \quad (4.6)$$

This uses only the first moment, which is the score from the conditional likelihood of s given x with h replaced by N_1/N . It is similar to Manski and McFaddens' (1981) conditional maximum likelihood estimator in the standard case-control or choice-based sampling set up. The asymptotic covariance matrix of this estimator is

$$V(\hat{\beta}_{CML}) = \Delta_{11}^{-1} - \Delta_{11}^{-1} \Delta_{12} [\bar{h}/q^2]^{-1} \Delta_{21} \Delta_{11}^{-1}. \quad (4.7)$$

This estimator is distributed independently of \hat{h} . Its inefficiency is revealed by comparison with (4.4) since $\Delta_{22} - \Delta_{21} \Delta_{11}^{-1} \Delta_{12}$ is non-negative definite.

$\sqrt{n}(\hat{h} - h)$ is distributed independently of $\hat{\beta}$ with variance \bar{h} .

5. THE LOGIT CASE

The logit model for P is of interest since it is widely used and there are known simplifications under this model in standard case-control sampling. The model is

$$P(x; \beta) = 1/(1 + \exp\{\beta_0 + \beta'_1 x\}).$$

Under standard case-control sampling the conditional probability of choice 1 given the covariate *and* the sampling scheme is

$$R_1(x; \beta) = 1/(1 + \exp\{\beta_0 + \log[q(1-h)/h(1-q)] + \beta'_1 x\})$$

which is the original logit model with intercept displaced. This is the reason why under standard case-control sampling with a logit model an investigator can proceed *as if* the data had been obtained by random sampling so far as inference about the covariate effects is concerned. But in the present application the conditional probability of stratum 1 given the covariate and the sampling scheme is

$$R_1(x; \beta) = 1/(1 + [q(1-h)/h] + \exp\{\beta_0 + \log[q(1-h)/h] + \beta'_1 x\}).$$

This is not a logit model. Thus it would be incorrect for an investigator to proceed to make inferences about covariate effects as if the data originated in random sampling.

Steinberg and Cardell (1991) have suggested an estimator for the logit model when q is known. They propose choosing β to maximize

$$L_{SC} = \sum_{n=1}^N (1 - s_n) \log(1 - P_n(\beta)) + \omega s_n \log[P_n(\beta)/(1 - P_n(\beta))]. \quad (5.1)$$

Here $\omega = q^*(1 - \hat{h})/\hat{h}$. In this section we shall give an interpretation of the Steinberg and Cardell (SC) estimator and comment on its properties.³ In the next section we report some Monte Carlo comparisons of this estimator and the efficient procedure.

³Steinberg and Cardell actually study a slightly different case where the population is finite, and the two samples, one containing observations with $y = 1$, and one randomly from the whole population, may partially overlap. The model we study can be viewed as a limit of their framework where the size of the population goes to infinity. They also gave a quite different justification for their estimator than the one which follows.

Consider a two stage estimation procedure. In the first stage a nonparametric estimate of the population joint distribution of choice and covariate is constructed. In the second stage, an estimate of β is formed by minimizing the Kullback-Leibler(KL) distance between the nonparametric estimate and a proposed parametric (logit) model. Let $u(y, x) = f(x)P(x)^y[1 - P(x)]^{1-y}$, the population joint distribution of choice and covariate. The second stage therefore minimizes

$$C = \sum_{y,x} \hat{u}(y, x) \log[\hat{u}(y, x)/u(y, x; \beta)] \quad (5.2)$$

where \hat{u} is the nonparametric estimate and $u(y, x; \beta)$ is the parametric model with a logit form for $P(x)$ depending on the parameter β . Dropping terms from (5.2) which do not involve β it may be written

$$\begin{aligned} C &= \sum_l \hat{P}_l \hat{f}_l \log P_l(\beta) + (1 - \hat{P}_l) \hat{f}_l \log(1 - P_l(\beta)) \\ &= \sum_l \hat{f}_l \log(1 - P_l(\beta)) + \hat{P}_l \hat{f}_l \log[P_l(\beta)/(1 - P_l(\beta))]. \end{aligned} \quad (5.3)$$

In this expression, $f_l = f(x^l)$, $P_l = P(x^l)$ and a caret indicates the nonparametric estimate.

Now consider nonparametric estimation of P and f . The log likelihood (3.3) with $g(x)$ multinomial leads us to such estimates. The ML estimate of $g(x)$ is $\hat{\pi}_l = n_l/N$. The nonparametric estimate of $R_1(x^l) = \hat{R}_{1l}$ is n_{1l}/n_l . Here n_l is the number of observations having covariate value x^l and n_{1l} is the number of observations having covariate x^l and originating from stratum 1. n_{0l} is similarly defined. Note that these estimates do satisfy the constraint (3.4) or (3.5) when $\hat{h} = N_1/N$ so they do in fact maximize the constrained log likelihood.

The definitions (3.2) and the definition of $\omega = q^*(1 - \hat{h})/\hat{h}$ enable us to go from estimates of R_1, g to estimates of f and Pf which are

$$\hat{P}_l \hat{f}_l = \frac{\omega n_{1l}}{N_0}; \quad \hat{f}_l = \frac{n_{0l}}{N_0}. \quad (5.4)$$

Note that the nonparametric estimator of $f(x)$ is the sample distribution from stratum 0, the random sample.

Inserting these estimates into the KL measure, (5.3), gives

$$\begin{aligned} C &= N_0^{-1} \sum_l n_{0l} \log(1 - P_l(\beta)) + \omega n_{1l} \log(P_l(\beta)/(1 - P_l(\beta))) \\ &= N_0^{-1} \sum_{n=1}^N (1 - s_n) \log(1 - P_n(\beta)) + \omega s_n \log[P_n(\beta)/(1 - P_n(\beta))] \end{aligned} \quad (5.5)$$

This is proportional to the Steinberg and Cardell criterion function, (5.1).

While the preceding argument is formally correct it suffers from the difficulty that the implicit 'nonparametric ML' estimate of P may lie outside the interval zero to one. This is obvious from the relation between R_1 and P given in (3.2) where, even though \hat{R}_1 is a proper probability there is no guarantee that \hat{P} is. This suggests that the Steinberg/Cardell estimator may behave poorly in small samples, even though when P is logit the criterion function (5.1) is globally concave.

It is interesting to look at the form of the Steinberg-Cardell estimator in more detail, as it explains some of the findings of the Monte Carlo study. Suppose that x is a scalar random variable, taking on two values, 0 and 1. Also, assume that β_0 is known. The first order condition for maximization of L_{SC} is

$$L_{SC}^{\beta} = \sum_{n=1}^N x_n [\omega s_n - (1 - s_n) P_n] = 0.$$

Since x is binary this becomes

$$L_{SC}^{\beta} = \omega S - P(N_1 - S) = 0, \quad (5.6)$$

where S is the number of the N_1 observations from stratum 1 having covariate value one and $P = 1/(1 + \exp\{\beta_0 + \beta_1\})$. Conditional on $x = 1$, S is Binomial ($N_1, R_1(1; \beta, q, h)$).

Equation (5.6) will have a finite solution for β_1 if and only if $\omega S \leq N_1 - S$, an event of probability less than one. As a particular example suppose that $\beta_1^* = 0$ and that equal numbers of observations come from each stratum, $h = 0.5$. Then, using the Normal approximation to the Binomial, we find

$$pr(\text{no finite solution for } \hat{\beta}_1) = 1 - \Phi\left(\sqrt{N_1} \frac{1 - q}{1 + q}\right)$$

Some values of this probability are given in the following table

‘Probabilities of No Solution ‘

N	N_1	q	Probability
100	50	0.90	0.355
400	200	0.90	0.228
1000	500	0.90	0.120
100	50	0.80	0.216
400	200	0.80	0.058
1000	500	0.80	0.0065

Under these circumstances the efficient GMM estimator can be expected to perform much better. The third moment compares q to the average value of $P(x; \beta)/(hP(x; \beta)/q + 1 - h)$. In this case with β_1 close to zero, this moment has very little variance, and gives an almost exact restriction on β_1 . This information is not used by the Steinberg-Cardell estimator. This is of course no proof that the GMM estimator will in fact perform better in practice. It relies on a first round of consistent estimates to get an estimate of the optimal weight matrix. The choice of the first round weight matrix does not matter asymptotically, but there is no guarantee that the first round estimator will actually converge. In practice however, we had no difficulty in obtaining convergence for the GMM estimator using prior knowledge of q .

In the Monte Carlo experiment x was chosen to have a bivariate normal distribution with zero means, unit variance and zero correlation. Three sets of parameter values were used: $(\beta_0, \beta_1, \beta_2)$ equal to $(0, 1, 1)$, $(0, 2, 0.5)$ and $(-1.89, 1, 1)$. The implied values for q were 0.5, 0.5 and 0.2. h was fixed at 0.5. The number of observations was in all simulations equal to 400. The number of replications was equal to 200 for each experiment. We report the averages of the 200 estimates (mean), the average of the asymptotic standard deviations (asd), the standard deviation of the 200 replications (ssd), the median, and the median of the absolute deviation from the median (mad). The results are reported in tables 1 to 3.

The SC estimator performed significantly worse than the efficient GMM estimator

Table 1: Design I									
$\beta_0 = 0.0, \beta_1 = 1.0, \beta_2 = 1.0, q = 0.5, h = 0.5$									
failure to converge	GMM(unknown q)			GMM(known q)			SC		
	11			0			9		
	β_0	β_1	β_2	β_0	β_1	β_2	β_0	β_1	β_2
mean	0.06	1.18	1.18	0.02	1.04	1.04	0.10	1.20	1.22
asd	33.98	14.43	11.41	0.10	0.26	0.26	0.42	0.77	0.77
ssd	0.98	0.48	0.49	0.10	0.29	0.27	0.34	0.64	0.61
med	0.06	1.09	1.11	0.01	1.01	1.02	0.04	1.05	1.10
mad	0.66	0.32	0.30	0.06	0.19	0.15	0.20	0.30	0.31

Table 2: Design II									
$\beta_0 = 0.0, \beta_1 = 2.0, \beta_2 = 0.5, q = 0.5, h = 0.5$									
failure to converge	GMM(unknown q)			GMM(known q)			SC		
	2			0			27		
	β_0	β_1	β_2	β_0	β_1	β_2	β_0	β_1	β_2
mean	-0.01	2.15	0.52	0.00	2.03	0.50	0.08	2.58	0.70
asd	17.00	16.35	5.02	0.13	0.38	0.25	0.92	5.15	2.10
ssd	0.81	0.66	0.31	0.13	0.38	0.26	0.46	1.98	0.93
med	-0.01	2.04	0.46	-0.01	1.98	0.49	0.00	2.10	0.51
mad	0.46	0.33	0.17	0.08	0.26	0.17	0.23	0.62	0.24

Table 3: Design III									
$\beta_0 = -1.89, \beta_1 = 1.0, \beta_2 = 1.000, q = 0.2, h = 0.5$									
failure to converge	GMM(unknown q)			GMM(known q)			SC		
	19			0			0		
	β_0	β_1	β_2	β_0	β_1	β_2	β_0	β_1	β_2
mean	-1.89	1.12	1.10	-1.87	1.04	1.03	-1.92	1.09	1.06
asd	133.68	61.76	61.90	0.09	0.18	0.18	0.20	0.32	0.32
ssd	0.75	0.31	0.26	0.10	0.20	0.18	0.20	0.36	0.36
med	-1.77	1.09	1.09	-1.87	1.04	1.03	-1.92	1.03	1.01
mad	0.39	0.16	0.18	0.06	0.12	0.13	0.12	0.20	0.17

proposed in this paper. In fact, in the first and third set of simulations 9 and 27 of the replications did not lead to convergence. The GMM estimator without knowledge of q did not converge for 11, 2 and 19 of the simulations. There were no problems with convergence of the GMM estimator with known q . The standard errors for the unknown q GMM estimator and the Steinberg-Cardell estimator reflect the convergence problems: they are markedly different from what one would expect given normality and given the median deviation from the mean. The finite sample properties of the known q GMM estimator seem satisfactory and reflect its theoretical asymptotic superiority to the Steinberg and Cardell estimator when the model is correctly specified.

6. SUMMARY AND CONCLUSIONS

We have given computationally simple and asymptotically efficient estimators in the contaminated sampling problem. When the marginal choice probability, q , is unknown the estimator maximizes a binary choice log likelihood and, if the covariate distribution is multinomial with known support, it is interpretable as a constrained maximum likelihood estimator. When the marginal choice probability is known the estimator solves a generalized method of moments problem. When the covariate distribution is multinomial with known support the estimator is asymptotically equivalent to a constrained maximum likelihood estimator. We also gave explicit forms for the asymptotic covariance matrices in both cases as well as for a conditional likelihood estimator applicable when q is known. Additional a priori information can be readily incorporated into the GMM procedure as long as it is expressible as a moment condition. Imbens and Lancaster(1992) gives further examples of this.

We have also discussed the logit model as a special case and compared numerically the properties of the estimators proposed in this paper with an alternative method suggested by Steinberg and Cardell(1991) which is applicable when q is known. When q is known the efficient generalized method of moments estimator exhibited satisfactory performance. The estimator of Steinberg and Cardell failed to exist in a significant fraction of simulations as did the efficient GMM procedure in the absence of knowledge of the marginal choice probability.

APPENDIX
PROOFS OF THEOREMS 1 AND 2

Consistency and asymptotic normality of the GMM estimators both when q is known and when it is unknown can be proved in a generalized method of moments framework as described by Hansen(1982) and Manski(1988). For instance, theorems 2.1 and 3.1 in Hansen(1982) prove consistency and asymptotic normality for generalized method of moments estimators. Conditions that ensure that the regularity conditions for these theorems are satisfied are: (i) compactness of sample and parameter spaces (with true parameters interior to the parameter space), (ii) continuity of $P(x; \beta)$ and its derivative with respect to β , (iii) uniqueness of the solution to $\mathcal{E}(\psi(\delta)) = 0$, and (iv) full rank of Δ and Γ .

The estimator of theorem 1 was derived initially for the case in which x has a discrete distribution with known, finite, support. The estimator was shown to be a maximum likelihood estimator in that case and therefore achieves the Cramer-Rao bound for regular estimators. This result can be extended to the continuous regressor case using the approach to semiparametric efficiency bounds of Begun, Hall, Huang and Wellner(1984).

From (3.1) the log density of a single observation is

$$\log g(s, x) = s \log P(x; \beta) - s \log q + s \log h + (1 - s) \log(1 - h) + \log f(x). \quad (\text{A1})$$

Consider a parametric submodel in which the unknown density $f(\cdot)$ is parametrized by η . In this submodel the scores for β and η are

$$S_\beta = s(p'_\beta/P - q_\beta/q); \quad S_\eta = -s(q_\eta/q) + f_\eta/f. \quad (\text{A2})$$

The tangent set, T , is of the form

$$d(x) - s(\mathcal{E}(d(x))/h)$$

where $d(x)$ is unrestricted apart from the requirement that $\int d(x) dF(x) = 0$. The efficient score is

$$S^* = (s - R(x; \beta))(p'_\beta/P + \delta/q),$$

where

$$\delta = -q \frac{\mathcal{E}(p'_\beta \bar{R}/P)}{\mathcal{E}(\bar{R})} = -\Delta_{12} \Delta_{22}^{-1}.$$

The inverse of the covariance matrix of S^* is the variance of the GMM estimator described in theorem 1.

The extension of this theorem to the case in which q is known is not yet available.

The claim in theorem 2 that the GMM estimator is asymptotically equivalent to the constrained ML estimator when the covariate distribution is discrete with known support is established by direct calculation using classical results on the covariance matrix of the constrained maximum likelihood estimator.

REFERENCES

- Begun, J. M., W. J. Hall, W-M. Huang and J. A. Wellner, (1983) 'Information and Asymptotic Efficiency in Parametric-Nonparametric models', *Annals of Statistics*, vol 11, 432-452.
- Breslow, N. E. and Day, N. (1980) 'Statistical Methods in Cancer Research, 1: The Analysis of Case-control Studies'. IARC, Lyon.
- Chamberlain, G., (1987), 'Asymptotic Efficiency in Estimation with Conditional Moment Restrictions', *Journal of Econometrics*, vol 34, 305-334, 1987.
- Cosslett, S. R., (1981b) 'Efficient Estimation of Discrete Choice Models', in C. F. Manski and D. McFadden, eds., *Structural Analysis of Discrete Data with Econometric Applications*, 51-111, MIT Press, Cambridge, MA,
- Cox, D., and D. Hinkley, (1974), *Theoretical Statistics*, Chapman and Hall, London
- Hansen, L. P., (1982), 'Large Sample Properties of Generalized Method of Moment Estimators', *Econometrica*, vol 50, 1029-1054.
- Heckman, J. J and Robb, (1984) 'Alternative Methods for Evaluating the Impact of Interventions', in J. J. Heckman and B. Singer, 'Longitudinal Analysis of Labor Market Data', Cambridge University Press.
- Hsieh, D. A., C. F. Manski and D. McFadden, (1985), 'Estimation of Response probabilities from Augmented Retrospective Observations', *Journal of the American Statistical Association* vol 80, 651-662,
- Imbens, G. W., (1990), 'An Efficient Method of Moments Estimator for Discrete Choice Models with Choice-based Sampling', discussion paper, CentER, Tilburg University, 1990
- Imbens, G. W., and T. Lancaster, (1991), 'Efficient Estimation and Stratified Sampling', Harvard Institute of Economic Research Discussion Paper.
- Imbens, G. W., and T. Lancaster, (1992), 'Combining Micro and Macro Data in Microeconomic Models', Harvard Institute of Economic Research Discussion Paper.
- Manski, C. F., (1988) *Analog Estimation Methods in Econometrics*, Chapman and Hall, New York, NY.

Manski, C. F., and S. R. Lerman, (1977), 'The Estimation of Choice Probabilities from Choice-based Samples', *Econometrica*, vol 45, 1977-1988,

Manski, C. F., and D. McFadden, (1981), 'Alternative Estimators and Sample Designs for Discrete Choice Analysis', in C. F. Manski and D. McFadden, eds., *Structural Analysis of Discrete Data with Econometric Applications*, 51-111, MIT Press, Cambridge, MA

Prentice, R. L., and R. Pyke, (1979), 'Logistic Disease Incidence Models and Case-Control Studies', *Biometrika*, 66, 3, 403-411.

Steinberg, D., and N. Cardell, (1991), 'Estimating Logistic Regression Models when the Dependent Variable has no Variance', forthcoming, *Communications in Statistics*.

Discussion Paper Series, CentER, Tilburg University, The Netherlands:

(For previous papers please consult previous discussion papers.)

No.	Author(s)	Title
9136	H. Bester and E. Petrakis	The Incentives for Cost Reduction in a Differentiated Industry
9137	L. Mirman, L. Samuelson and E. Schlee	Strategic Information Manipulation in Duopolies
9138	C. Dang	The D'_2 -Triangulation for Continuous Deformation Algorithms to Compute Solutions of Nonlinear Equations
9139	A. de Zeeuw	Comment on "Nash and Stackelberg Solutions in a Differential Game Model of Capitalism"
9140	B. Lockwood	Border Controls and Tax Competition in a Customs Union
9141	C. Fershtman and A. de Zeeuw	Capital Accumulation and Entry Deterrence: A Clarifying Note
9142	J.D. Angrist and G.W. Imbens	Sources of Identifying Information in Evaluation Models
9143	A.K. Bera and A. Ullah	Rao's Score Test in Econometrics
9144	B. Melenberg and A. van Soest	Parametric and Semi-Parametric Modelling of Vacation Expenditures
9145	G. Imbens and T. Lancaster	Efficient Estimation and Stratified Sampling
9146	Th. van de Klundert and S. Smulders	Reconstructing Growth Theory: A Survey
9147	J. Greenberg	On the Sensitivity of Von Neuman and Morgenstern Abstract Stable Sets: The Stable and the Individual Stable Bargaining Set
9148	S. van Wijnbergen	Trade Reform, Policy Uncertainty and the Current Account: A Non-Expected Utility Approach
9149	S. van Wijnbergen	Intertemporal Speculation, Shortages and the Political Economy of Price Reform
9150	G. Koop and M.F.J. Steel	A Decision Theoretic Analysis of the Unit Root Hypothesis Using Mixtures of Elliptical Models
9151	A.P. Barten	Consumer Allocation Models: Choice of Functional Form
9152	R.T. Baillie, T. Bollerslev and M.R. Redfearn	Bear Squeezes, Volatility Spillovers and Speculative Attacks in the Hyperinflation 1920s Foreign Exchange

No.	Author(s)	Title
9153	M.F.J. Steel	Bayesian Inference in Time Series
9154	A.K. Bera and S. Lee	Information Matrix Test, Parameter Heterogeneity and ARCH: A Synthesis
9155	F. de Jong	A Univariate Analysis of EMS Exchange Rates Using a Target
9156	B. le Blanc	Economies in Transition
9157	A.J.J. Talman	Intersection Theorems on the Unit Simplex and the Simplotope
9158	H. Bester	A Model of Price Advertising and Sales
9159	A. Özcam, G. Judge, A. Bera and T. Yancey	The Risk Properties of a Pre-Test Estimator for Zellner's Seemingly Unrelated Regression Model
9160	R.M.W.J. Beetsma	Bands and Statistical Properties of EMS Exchange Rates: A Monte Carlo Investigation of Three Target Zone Models Zone Model
9161	A.M. Lejour and H.A.A. Verbon	Centralized and Decentralized Decision Making on Social Insurance in an Integrated Market Multilateral Institutions
9162	S. Bhattacharya	Sovereign Debt, Creditor-Country Governments, and
9163	H. Bester, A. de Palma, W. Leininger, E.-L. von Thadden and J. Thomas	The Missing Equilibria in Hotelling's Location Game
9164	J. Greenberg	The Stable Value
9165	Q.H. Vuong and W. Wang	Selecting Estimated Models Using Chi-Square Statistics
9166	D.O. Stahl II	Evolution of Smart _n Players
9167	D.O. Stahl II	Strategic Advertising and Pricing with Sequential Buyer Search
9168	T.E. Nijman and F.C. Palm	Recent Developments in Modeling Volatility in Financial Data
9169	G. Asheim	Individual and Collective Time Consistency
9170	H. Carlsson and E. van Damme	Equilibrium Selection in Stag Hunt Games
9201	M. Verbeek and Th. Nijman	Minimum MSE Estimation of a Regression Model with Fixed Effects from a Series of Cross Sections
9202	E. Bomhoff	Monetary Policy and Inflation
9203	J. Quiggin and P. Wakker	The Axiomatic Basis of Anticipated Utility; A Clarification

No.	Author(s)	Title
9204	Th. van de Klundert and S. Smulders	Strategies for Growth in a Macroeconomic Setting
9205	E. Siandra	Money and Specialization in Production
9206	W. Härdle	Applied Nonparametric Models
9207	M. Verbeek and Th. Nijman	Incomplete Panels and Selection Bias: A Survey
9208	W. Härdle and A.B. Tsybakov	How Sensitive Are Average Derivatives?
9209	S. Albæk and P.B. Overgaard	Upstream Pricing and Advertising Signal Downstream Demand
9210	M. Cripps and J. Thomas	Reputation and Commitment in Two-Person Repeated Games
9211	S. Albæk	Endogenous Timing in a Game with Incomplete Information
9212	T.J.A. Storcken and P.H.M. Ruys	Extensions of Choice Behaviour
9213	R.M.W.J. Beetsma and F. van der Ploeg	Exchange Rate Bands and Optimal Monetary Accommodation under a Dirty Float
9214	A. van Soest	Discrete Choice Models of Family Labour Supply
9215	W. Güth and K. Ritzberger	On Durable Goods Monopolies and the (Anti-) Coase-Conjecture
9216	A. Simonovits	Indexation of Pensions in Hungary: A Simple Cohort Model
9217	J.-L. Ferreira, I. Gilboa and M. Maschler	Credible Equilibria in Games with Utilities Changing during the Play
9218	P. Borm, H. Keiding, R. Mclean, S. Oortwijn and S. Tijs	The Compromise Value for NTU-Games
9219	J.L. Horowitz and W. Härdle	Testing a Parametric Model against a Semiparametric Alternative
9220	A.L. Bovenberg	Investment-Promoting Policies in Open Economies: The Importance of Intergenerational and International Distributional Effects
9221	S. Smulders and Th. van de Klundert	Monopolistic Competition, Product Variety and Growth: Chamberlin vs. Schumpeter
9222	H. Bester and E. Petrakis	Price Competition and Advertising in Oligopoly

No.	Author(s)	Title
9223	A. van den Nouweland, M. Maschler and S. Tijs	Monotonic Games are Spanning Network Games
9224	H. Suehiro	A "Mistaken Theories" Refinement
9225	H. Suehiro	Robust Selection of Equilibria
9226	D. Friedman	Economically Applicable Evolutionary Games
9227	E. Bomhoff	Four Econometric Fashions and the Kalman Filter Alternative - A Simulation Study
9228	P. Borm, G.-J. Otten and H. Peters	Core Implementation in Modified Strong and Coalition Proof Nash Equilibria
9229	H.G. Bloemen and A. Kapteyn	The Joint Estimation of a Non-Linear Labour Supply Function and a Wage Equation Using Simulated Response Probabilities
9230	R. Beetsma and F. van der Ploeg	Does Inequality Cause Inflation? - The Political Economy of Inflation, Taxation and Government Debt
9231	G. Almekinders and S. Eijffinger	Daily Bundesbank and Federal Reserve Interventions - Do they Affect the Level and Unexpected Volatility of the DM/\$-Rate?
9232	F. Vella and M. Verbeek	Estimating the Impact of Endogenous Union Choice on Wages Using Panel Data
9233	P. de Bijl and S. Goyal	Technological Change in Markets with Network Externalities
9234	J. Angrist and G. Imbens	Average Causal Response with Variable Treatment Intensity
9235	L. Meijdam, M. van de Ven and H. Verbon	Strategic Decision Making and the Dynamics of Government Debt
9236	H. Houba and A. de Zeeuw	Strategic Bargaining for the Control of a Dynamic System in State-Space Form
9237	A. Cameron and P. Trivedi	Tests of Independence in Parametric Models: With Applications and Illustrations
9238	J.-S. Pischke	Individual Income, Incomplete Information, and Aggregate Consumption
9239	H. Bloemen	A Model of Labour Supply with Job Offer Restrictions
9240	F. Drost and Th. Nijman	Temporal Aggregation of GARCH Processes
9241	R. Gilles, P. Ruys and J. Shou	Coalition Formation in Large Network Economies
9242	P. Kort	The Effects of Marketable Pollution Permits on the Firm's Optimal Investment Policies

No.	Author(s)	Title
9243	A.L. Bovenberg and F. van der Ploeg	Environmental Policy, Public Finance and the Labour Market in a Second-Best World
9244	W.G. Gale and J.K. Scholz	IRAs and Household Saving
9245	A. Bera and P. Ng	Robust Tests for Heteroskedasticity and Autocorrelation Using Score Function
9246	R.T. Baillie, C.F. Chung and M.A. Tieslau	The Long Memory and Variability of Inflation: A Reappraisal of the Friedman Hypothesis
9247	M.A. Tieslau, P. Schmidt and R.T. Baillie	A Generalized Method of Moments Estimator for Long-Memory Processes
9248	K. Wärneryd	Partisanship as Information
9249	H. Huizinga	The Welfare Effects of Individual Retirement Accounts
9250	H.G. Bloemen	Job Search Theory, Labour Supply and Unemployment Duration
9251	S. Eijffinger and E. Schaling	Central Bank Independence: Searching for the Philosophers' Stone
9252	A.L. Bovenberg and R.A. de Mooij	Environmental Taxation and Labor-Market Distortions
9253	A. Lusardi	Permanent Income, Current Income and Consumption: Evidence from Panel Data
9254	R. Beetsma	Imperfect Credibility of the Band and Risk Premia in the European Monetary System
9301	N. Kahana and S. Nitzan	Credibility and Duration of Political Contests and the Extent of Rent Dissipation
9302	W. Güth and S. Nitzan	Are Moral Objections to Free Riding Evolutionarily Stable?
9303	D. Karotkin and S. Nitzan	Some Peculiarities of Group Decision Making in Teams
9304	A. Lusardi	Euler Equations in Micro Data: Merging Data from Two Samples
9305	W. Güth	A Simple Justification of Quantity Competition and the Cournot-Oligopoly Solution
9306	B. Peleg and S. Tijs	The Consistency Principle For Games In Strategic Form
9307	G. Imbens and T. Lancaster	Case Control Studies With Contaminated Controls

P.O. BOX 90153, 5000 LE TILBURG, THE NETHERLANDS

Bibliotheek K. U. Brabant



17 000 01117405 0