

# LARGE SAMPLE THEORY OF MAXIMUM LIKELIHOOD ESTIMATES IN SEMIPARAMETRIC BIASED SAMPLING MODELS<sup>1</sup>

BY PETER B. GILBERT

*Harvard School of Public Health*

Vardi [*Ann. Statist.* **13** 178–203 (1985)] introduced an  $s$ -sample biased sampling model with known selection weight functions, gave a condition under which the common underlying probability distribution  $G$  is uniquely estimable and developed simple procedure for computing the nonparametric maximum likelihood estimator (NPMLE)  $\mathbb{G}_n$  of  $G$ . Gill, Vardi and Wellner thoroughly described the large sample properties of Vardi's NPMLE, giving results on uniform consistency, convergence of  $\sqrt{n} (\mathbb{G}_n - G)$  to a Gaussian process and asymptotic efficiency of  $\mathbb{G}_n$ . Gilbert, Lele and Vardi considered the class of semiparametric  $s$ -sample biased sampling models formed by allowing the weight functions to depend on an unknown finite-dimensional parameter  $\theta$ . They extended Vardi's estimation approach by developing a simple two-step estimation procedure in which  $\hat{\theta}_n$  is obtained by maximizing a profile partial likelihood and  $\mathbb{G}_n \equiv \mathbb{G}_n(\hat{\theta}_n)$  is obtained by evaluating Vardi's NPMLE at  $\hat{\theta}_n$ . Here we examine the large sample behavior of the resulting joint MLE  $(\hat{\theta}_n, \mathbb{G}_n)$ , characterizing conditions on the selection weight functions and data in order that  $(\hat{\theta}_n, \mathbb{G}_n)$  is uniformly consistent, asymptotically Gaussian and efficient.

Examples illustrated here include clinical trials (especially HIV vaccine efficacy trials), choice-based sampling in econometrics and case-control studies in biostatistics.

**1. Introduction: semiparametric biased sampling models and the MLE.** Vardi (1985) developed methodology for the  $s$ -sample biased sampling model with known selection bias weight functions. Gilbert, Lele and Vardi (1999) considered an extension of this model which allows the weight functions to depend on an unknown finite-dimensional parameter  $\theta$ . This model comprises three components: a probability measure  $G$  defined on a sample space  $\mathbf{Y}$  with  $\sigma$ -field of subsets  $\mathbf{B}$ , a set of nonnegative (measurable) "stratum" weight functions  $w_1, \dots, w_s$  defined on  $\mathbf{Y} \times \Theta$  and selection probabilities  $\lambda_i, i = 1, \dots, s$ , with  $\sum_{i=1}^s \lambda_i = 1$ . The parameter  $\theta$  is common to all  $s$  weight functions  $w_i(\cdot, \theta)$ , which are assumed to be nonnegative and of a known parametric form. Let  $g \equiv dG/d\mu$  for some measure  $\mu$  dominating  $G$ . The data are assumed to be an i.i.d. sample  $X_k = (I_k, Y_k), k = 1, \dots, n$ , from the semiparametric biased sampling model  $\mathbf{P} \equiv \{P_{(\theta, G)}; \theta \in \Theta, G \in \mathbf{G}\}$ , defined on

---

Received January 1998; revised June 1999.

<sup>1</sup>Research supported by NIH Grants 5-U01-AI38855, 5-U01-AI28076 and AI24643.

AMS 1991 subject classifications. Primary 60G05, 62F05; secondary 62G20, 62G30.

*Key words and phrases.* Asymptotic theory, choice-based sampling, clinical trials, empirical processes, generalized logistic regression, HIV vaccine trial, nonparametric maximum likelihood, selection bias models, Vardi's estimator.

$\mathbf{X} \equiv \{1, \dots, s\} \times \mathbf{Y}$ , with density given by

$$(1.1) \quad p(x, \theta, G) = p(i, y, \theta, G) = \lambda_i \frac{w_i(y, \theta)}{W_i(\theta, G)} g(y).$$

Here  $W_i(\theta, G)$  is the  $i$ th normalizing function given by

$$(1.2) \quad W_i(\theta, G) \equiv \int_{\mathbf{Y}} w_i(u, \theta) dG(u),$$

assumed to be positive and finite for all  $\theta$  in the parameter space  $\Theta$ . The random variable  $I \in \{1, \dots, s\}$  denotes the stratum, selected with probability  $\lambda_i$ , and, conditional on  $I = i$ , the probability measure  $F_i$  of  $Y$  under the biased sampling model satisfies

$$(1.3) \quad F_i(A, \theta, G) \equiv W_i^{-1}(\theta, G) \int_A w_i(u, \theta) dG(u), \quad A \in \mathbf{B}, \quad i = 1, \dots, s.$$

In Vardi's (1985) original treatment, the weight functions were assumed completely known (independent of  $\theta$ ), but there are many practical situations in which a complete specification of the  $w_i$ 's is too restrictive, but the weight functions can realistically be assumed to belong to a parametric family. We illustrate with three examples.

**EXAMPLE 1.1** (Univariate generalized logistic regression model). As described in the introduction of Gilbert, Lele and Vardi (1999), the generalized logistic regression (GLR) model is useful for assessing from viral data taken from breakthrough infections in preventive human immunodeficiency virus (HIV) vaccine efficacy trials how vaccine protection against infection varies by characteristics of challenge HIV. Suppose viruses are ordered by some "distance," quantifiable in a variety of ways, for instance by the percent nucleotide or amino acid mismatch in a particular gene. Let  $Y \in [0, \infty)$  be a random variable denoting the distance between a virus that has been isolated from an infected trial participant and the virus strain represented in the vaccine. The GLR model is an  $s$ -sample semiparametric biased sampling model, where  $s$  is the number of intervention arms,  $G$  is the distribution of  $Y$  in the baseline placebo group (group  $s$ , say) and  $\theta$  describes how strongly vaccine protection depends on strain distance. For instance, in the two-sample GLR model, the groups are vaccine and placebo. Letting  $F_v$  denote the distribution of  $Y$  in the vaccine group, the GLR model relates  $G$  and  $F_v$  in the following way:

$$(1.4) \quad F_v(y) = \frac{\int_0^y \exp\{h(u, \theta)\} dG(u)}{\int_0^\infty \exp\{h(u, \theta)\} dG(u)}, \quad y \in [0, \infty),$$

where  $h(y, \theta)$  is a given function (most simply taken to be linear,  $h(y, \theta) = y\theta$ ) and  $\theta$  is an unknown  $d$ -dimensional "differential vaccine efficacy" parameter. This is a biased sampling model (1.1)–(1.3), with  $\mathbf{Y} = [0, \infty)$ ,  $F_1 = F_v$ ,  $F_2 = G$ ,  $w_1(y, \theta) = \exp\{h(y, \theta)\}$ ,  $w_2(y, \theta) = 1$  and  $\lambda_i \equiv n_i/n$ ,  $i = 1, 2$ , the sampling fraction of infections in the trial from group  $i$ . After reparametrization, the

ordinary logistic regression model with case-control sampling of covariates  $y$  can be written in the form of a two-sample GLR model with  $h(y, \theta) = \theta_1 + \theta_2^T y$ , where  $G$  is the distribution of  $y$  for controls.

The general  $s$ -sample GLR model is formed by supposing that observations from the  $i$ th sample are distributed i.i.d.  $F_i$  with

$$(1.5) \quad F_i(y) = \frac{\int_0^y \exp\{h_i(u, \theta)\} dG(u)}{\int_0^\infty \exp\{h_i(u, \theta)\} dG(u)}, \quad y \in [0, \infty),$$

where  $G$  is the (baseline) distribution for the  $s$ th sample and  $h_i(y, \theta)$ ,  $i = 1, \dots, s$ , are given functions. For identifiability and unique estimability, we take  $h_s \equiv 0$ , so that  $F_s = G$ . A practical choice for the functions  $\{h_i\}$  is  $h_i(y, \theta) = \sum_{k=1}^d h_{ik}(y)\theta_k$ , where the  $h_{ik}$  are given functions of  $y$  independent of  $\theta$ .

The GLR model can be used to assess how protection of a vaccine against any heterogeneous pathogen (not just HIV) depends on features of the pathogen, and more generally for investigating how the relative efficacy of two treatments varies by some random variable measured only on disease cases. Among many applications are the study of the relationship between a measure of genetic evolution of HIV in an infected person and the relative efficacy of two antiretroviral regimens, and the study of how the efficacy of a new treatment relative to a control treatment changes over time. The broad utility of the GLR model, especially for analyzing clinical trial data, is discussed more extensively in Gilbert, Self and Ashby (1998) and Gilbert, Lele and Vardi (1999).

**EXAMPLE 1.2** (Multivariate generalized logistic regression model). Since HIV vaccine protection may vary with variations in several viral features, it is of practical interest to assess differential vaccine efficacy by modeling a multidimensional distance  $Y \equiv (Y_1, \dots, Y_k)^T$  according to a GLR model. This is easily done in the same manner as the univariate model in Example 1.1. For instance, a simple two-sample GLR model with a bivariate distance is specified by

$$(1.6) \quad F_v(y_1, y_2) = \frac{\int_0^{y_1} \int_0^{y_2} \exp\{u_1\theta_1 + u_2\theta_2 + u_1u_2\theta_3\} dG(u_1, u_2)}{\int_0^\infty \int_0^\infty \exp\{u_1\theta_1 + u_2\theta_2 + u_1u_2\theta_3\} dG(u_1, u_2)}$$

for  $y_1, y_2 \in [0, \infty)$ . This is a two-sample biased sampling model with  $\mathbf{Y} = [0, \infty) \times [0, \infty)$ ,  $d = 3$ ,  $w_1(y, \theta) = \exp\{y_1\theta_1 + y_2\theta_2 + y_1y_2\theta_3\}$  and  $w_2(y, \theta) \equiv 1$ . The interaction parameter  $\theta_3$  will be zero if and only if the sum of the marginal vaccine effects against a strain with  $Y_1 = y_1$  and a strain with  $Y_2 = y_2$  equals the joint vaccine effect against a strain with  $Y_1 = y_1$  and  $Y_2 = y_2$ . Setting  $y_1 = 0$  ( $y_2 = 0$ ) recovers the marginal univariate GLR model for  $Y_2$  ( $Y_1$ ).

The full  $k$ -variate GLR model applies to random variables defined on  $\mathbf{Y} = [0, \infty)^k$ , and can be laid out similarly to (1.6), with  $\binom{k}{2}$  two-way interaction

parameters,  $\binom{k}{3}$  three-way interaction parameters and so on. Parsimonious versions can be considered by setting various interaction parameters to zero.

EXAMPLE 1.3 (Choice-based sampling in econometrics; case-control studies in biostatistics). This example is adapted from Example 4.4 of Gill, Vardi and Wellner (1988). Suppose that  $X = (J, Y)$ , where  $J$  is discrete with values on  $\{1, \dots, M\}$  and  $Y \sim G$  with density  $g$  with respect to  $\mu$  is a covariate vector with values in  $\mathbf{Y} \subset$  some  $R^d$ . The prospective, unbiased model  $H$  has density  $h(j, y) = p_\theta(j|y)g(y)$ , so that

$$H(\{j\} \times A) = \int_A p_\theta(j|y) dG(y)$$

for  $j = 1, \dots, M$  and  $A \in \mathbf{B}(R^d)$ , where  $p_\theta(j|y) = P_\theta(J = j|Y = y)$  is a parametric model. The multinomial logistic regression model is a common choice of  $p_\theta$ , given by

$$(1.7) \quad p_\theta(j|y) = \frac{\exp(\alpha_j + \beta_j^T y)}{\sum_{j'=1}^M \exp(\alpha_{j'} + \beta_{j'}^T y)},$$

with  $\theta = (\alpha, \beta) \in R^{(d+1)M}$ ,  $\alpha_M = \beta_M = 0$ .

The retrospective, biased sampling model  $F$  is obtained from  $H$  by the biasing functions  $w_i(y) = 1_{D_i}(y)$ , where  $D_i \subset \{1, \dots, M\}$  for  $i = 1, \dots, s$ . For instance, consider the “pure choice-based sampling model” as described in Cosslett (1981), where the strata  $D_i$  are defined by  $D_i \equiv \{i\}$ ,  $i = 1, \dots, s \equiv M$ . As in Gill, Vardi and Wellner [(1988), page 1094], Vardi’s (1985) necessary and sufficient condition for identifiability of  $G$  fails. Manski and Lerman (1977) avoid this problem by assuming that the “aggregate shares”

$$(1.8) \quad H(j, \theta, G) \equiv H(\{j\} \times \mathbf{Y}, \theta, G) = \int_{\mathbf{Y}} p_\theta(j|y) dG(y), \quad j = 1, \dots, M,$$

are known. For this biasing system, one can view  $F$  as a biased distribution derived from  $G$  with new weight functions  $w_j^*(y, \theta) = p_\theta(j|y)$ ,  $j = 1, \dots, M$ . Then  $G$  in model (1.1)–(1.3) is usually identifiable, as all that is required is  $\int 1_{[p_\theta(j|y) > 0]} p_\theta(j|y) dG(y) > 0$ ,  $j = 1, \dots, M$  (see Proposition 1.1 of Gill, Vardi and Wellner (1988)). If  $\theta$  is known, the methods of Vardi (1985) and Gill, Vardi and Wellner (1988) apply. However, in many practical applications  $\theta$  will be unknown. Since all  $M$  weight functions depend on  $\theta$ , the resulting model is not a member of the class of biased sampling models primarily studied here, in which one weight function is independent of  $\theta$ . However, the extra assumption that the normalizing constants  $H(j, \theta, G)$  are known implies that the estimation procedure considered here still produces the maximum likelihood estimator (MLE) of  $(\theta, G)$ , with limiting behavior characterized by results given here.

The pure choice-based sampling design is also frequently used in case-control studies in biostatistics, where the  $j$ ’s often represent different disease categories. In the biostatistics application, the aim is to estimate the odds ratios  $\{\beta_j\}$  of (1.7). Given known  $H(j, \theta, G)$  as in (1.8), the asymptotic behavior

of these odds ratios and  $G$  are described by results developed here. The results also apply to choice-based models with general  $D_i$ 's for which the biased distribution  $F$  has density

$$(1.9) \quad \begin{aligned} f(i, j, y) &= \lambda_i \frac{\mathbf{1}_{D_i}(j) p_\theta(j|y) g(y)}{\int \sum_{j'=1}^M \mathbf{1}_{D_i}(j') p_\theta(j'|y') g(y') d\mu(y')} \\ &\equiv \lambda_i \frac{w_i^{**}(j, y, \theta) g(y)}{\int \sum_{j'=1}^M w_i^{**}(j', y', \theta) dG(y')}. \end{aligned}$$

Other examples and applications of semiparametric biased sampling models are described in Gilbert (1996) and Gilbert, Lele and Vardi (1999). Before presenting methodology for estimation of  $\theta$  and  $G$  in these models, we note that Qin (1998) considered the closely related semiparametric two-sample density ratio model. Our approach resembles that taken by Qin (1998) in its utilization of a profile likelihood, but substantially differs in that it explicitly extends the original approach of Vardi (1985) and Gill, Vardi and Wellner (1988), thereby including estimation of the normalizing constants as a central component. In addition, the treatment here is more general in that  $s$  samples rather than two are handled, and arbitrary selection weight functions are allowed.

The first matter at hand is identifiability. As Gill, Vardi and Wellner (1988) is repeatedly referenced, henceforth it is abbreviated GVW.

**1.1. Identifiability.** In the case of known  $\theta$ , GVW characterized necessary and sufficient conditions for  $G$  to be identifiable, namely, that the sample space  $\mathbf{Y}$  equals  $[y: w_i(y, \theta) > 0 \text{ for some } i = 1, \dots, s]$  and a certain graph is connected. For fixed  $\theta \in \Theta$ , define a graph  $\mathbf{G}^*(\theta)$  on the  $s$  vertices  $i = 1, \dots, s$  by identifying vertex  $i$  with  $k$  if and only if

$$\int \mathbf{1}_{[w_i(y, \theta) > 0]} \mathbf{1}_{[w_k(y, \theta) > 0]} dG(y) > 0.$$

The graph  $\mathbf{G}^*(\theta)$  is *connected* if every pair of vertices is connected by a path. In what follows, assume these conditions hold for all  $\theta$  in some neighborhood  $\Theta_0 \subset \Theta$  of the true  $\theta_0$ . Notice that these two conditions hold automatically if all of the weight functions are strictly positive.

Now consider identifiability of unknown  $\theta$  and  $G$ . When there is only one sample, the model is rarely identifiable. To our knowledge, the only class of one-sample models (with no restrictions on  $G$ ) known to be identifiable are those with weight function whose domain depends on  $\theta$  [see Gilbert, Lele and Vardi (1999), Theorem 1]. For the  $s$ -sample model, with  $s \geq 2$ , a large class of models is identifiable. When one of the weight functions is independent of  $\theta$ , the class of identifiable models can be characterized by a simple condition on the weight functions. The following result is proved in Gilbert, Lele and Vardi (1999).

**THEOREM 1.1.** *Let  $s \geq 2$ , with  $w_s$  independent of  $\theta$ . Suppose the sample space  $\mathbf{Y}$  equals  $\cup_{i=1}^s [y: w_i(y, \theta) > 0 \text{ for all } \theta \in \Theta]$  and  $\mathbf{G}^*(\theta)$  is connected for all  $\theta$  in some neighborhood  $\Theta_0 \subset \Theta$  of the true  $\theta_0$ . Then the  $s$ -sample biased sampling model is identifiable if and only if the following condition holds: for all  $\tilde{\theta}, \theta \in \Theta_0$  with  $\tilde{\theta} \neq \theta$ , there is at least one weight function  $w_i, i \in \{1, \dots, s-1\}$ , such that  $w_i(y, \tilde{\theta})$  and  $w_i(y, \theta)$  are linearly independent as functions of  $y$ .*

Gilbert, Lele and Vardi (1999) describe a large class of models that satisfy the conditions of Theorem 1.1. For instance, all interesting GLR models are identifiable by this theorem.

Henceforth, assume the identifiability conditions of Theorem 1.1, so that we consider identifiable biased sampling models (1.1)–(1.3) with multiple samples, in which at least one of the weight functions is independent of  $\theta$ . In fact, throughout we require that the  $s$ th weight function is constant, which is used to ensure unique estimability by the employed estimation procedure.

**1.2. Maximum likelihood estimation.** The MLE that we study is computed by a procedure introduced by Gilbert, Lele and Vardi (1999). We sketch it here. Following their notation, denote the size of the  $i$ th sample by  $n_i$ , the total sample size by  $n = \sum_{i=1}^s n_i$ , and the  $i$ th sampling fraction by  $\lambda_{ni} = n_i/n$ . Let  $t_1, \dots, t_h$  be the distinct observed  $Y$  values, with multiplicities  $r_1, \dots, r_h$ . Let  $n_{ij}, i = 1, \dots, s, j = 1, \dots, h$ , be the number of observations from the  $i$ th group with value  $t_j$ . Notice that  $h$  is random. Then the likelihood of the data observed according to (1.1)–(1.3) is

$$(1.10) \quad L_n(\theta, G|\underline{x}) = \prod_{i=1}^s \prod_{j=1}^h \left[ \frac{w_i(t_j, \theta)G\{t_j\}}{W_i(\theta, G)} \right]^{n_{ij}},$$

where  $\underline{x}$  represents dependency on the observed data  $\{X_k = (I_k, Y_k), k = 1, \dots, n\}$ , equivalent to  $\{t_j, n_{ij}: i = 1, \dots, s, j = 1, \dots, h\}$ . Define  $w_{ij}(\theta) = w_i(t_j, \theta)$ . Consider a partial likelihood defined by

$$(1.11) \quad L_{n1}(\theta, \underline{V}|\underline{x}) = \prod_{i=1}^s \prod_{j=1}^h \left[ \frac{w_{ij}(\theta)V_i^{-1}}{\sum_{k=1}^s \lambda_{nk}w_{kj}(\theta)V_k^{-1}} \right]^{n_{ij}},$$

with  $\underline{V} = (V_1, \dots, V_s)^T$  given by  $V_i = W_i(\theta, G)/W_s(\theta, G), i = 1, \dots, s$ , where the dependency of  $V_i$  on  $\theta$  and  $G$  is suppressed. Notice that  $L_{n1}$  depends on  $G$  only through the normalizing constants. When  $w_s$  is a constant function (or when  $w_s$  depends on  $\theta$  but  $W_s(\theta, G)$  is a known constant), the following maximum partial likelihood estimation procedure yields the estimate  $(\hat{\theta}_n, \mathbb{G}_n)$  which maximizes the full likelihood:

- 1 Maximize  $L_{n1}$  over  $\theta$  and  $\underline{V}$ , subject to  $V_1 > 0, V_2 > 0, \dots, V_{s-1} > 0, V_s = 1$  to obtain  $(\hat{\theta}_n, \underline{V}_n)$ .
- 2 Compute Vardi's NPML  $\mathbb{G}_n \equiv \mathbb{G}_n(\hat{\theta}_n)$  from data with "known" weight functions  $w_i(\cdot, \hat{\theta}_n)$ .

3 Estimate  $W_i$  by  $\mathbb{W}_{ni} \equiv \mathbb{W}_{ni}(\hat{\theta}_n) = \int w_i(u, \hat{\theta}_n) d\mathbb{G}_n(u)$ ,  $i = 1, \dots, s$ .

The key assumption for this procedure to produce the unique maximum is that a graphical condition holds, defined in Vardi (1985) and restated here. For fixed, known  $\theta \in \Theta$ , consider the graph  $\mathbf{G}(\theta)$  on the  $s$  vertices  $i = 1, \dots, s$  defined as follows. Define a directed edge from a vertex  $i$  to a vertex  $k$ ,  $i \leftrightarrow k$ , if and only if  $\sum_{j=1}^h w_{ij}(\theta)n_{kj} > 0$ . The graph  $\mathbf{G}(\theta)$  is *strongly connected* if, for every pair  $(i, k)$ , there exists a directed path from  $i$  to  $k$  and a directed path from  $k$  to  $i$ .

Theorem 1.2, proved in Gilbert, Lele and Vardi (1999), asserts that the above estimation procedure yields the MLE of  $(\theta, G)$ .

**THEOREM 1.2.** *Suppose  $s \geq 2$  with  $w_s$  a constant function, and the identifiability conditions of Theorem 1.1 hold. Further suppose the graph  $\mathbf{G}(\theta)$  is strongly connected for all  $\theta$  in some neighborhood  $\Theta_0 \subset \Theta$  of  $\theta_0$ , and that  $(\hat{\theta}_n, \underline{\mathbb{V}}_n, \mathbb{G}_n)$  is obtained from procedure 1–3. If  $(\hat{\theta}_n, \underline{\mathbb{V}}_n)$  uniquely maximizes the partial likelihood (1.11), then  $(\hat{\theta}_n, \mathbb{G}_n)$  uniquely maximizes the full likelihood  $L_n(\theta, G|\underline{\mathbf{x}})$  of (1.10).*

Procedure 1–3 can be carried out as follows. Step 1 can be accomplished via profile likelihood. For fixed  $\theta \in \Theta$ , let  $\underline{\mathbb{V}}_n(\theta) = (\mathbb{V}_{n1}(\theta), \dots, \mathbb{V}_{ns-1}(\theta), 1)^T$  be the unique solution of

$$(1.12) \quad \mathbb{H}_{ni}(V_1(\theta), \dots, V_s(\theta)) \equiv V_i^{-1}(\theta) \sum_{j=1}^h \frac{r_j w_{ij}(\theta)}{\sum_{k=1}^s n_k w_{kj}(\theta) V_k^{-1}(\theta)} = 1,$$

$i = 1, \dots, s-1$ , in the region  $V_1(\theta) > 0, \dots, V_{s-1}(\theta) > 0, V_s = 1$ . Vardi (1985) proved that (1.12) has a unique solution if and only if the graph  $\mathbf{G}(\theta)$  is strongly connected. The estimator  $\hat{\theta}_n$  is the argument which maximizes the profile partial likelihood  $L_{n1\text{pro}}$  defined by

$$(1.13) \quad \begin{aligned} L_{n1\text{pro}}(\theta|\underline{\mathbf{x}}) &= L_{n1\text{pro}}(\theta, \underline{\mathbb{V}}_n(\theta)|\underline{\mathbf{x}}) \\ &= \prod_{i=1}^s \prod_{j=1}^h \left[ \frac{w_{ij}(\theta) \mathbb{V}_{ni}^{-1}(\theta)}{\sum_{k=1}^s \lambda_{nk} w_{kj}(\theta) \mathbb{V}_{nk}^{-1}(\theta)} \right]^{n_{ij}}. \end{aligned}$$

Set  $\underline{\mathbb{V}}_n = \underline{\mathbb{V}}_n(\hat{\theta}_n)$ . Then Step 2 proceeds by setting  $\hat{p} = p(\hat{\theta}_n)$ , where

$$p_j(\theta) \propto r_j \left/ \sum_{k=1}^s n_k w_{kj}(\theta) \mathbb{V}_{nk}^{-1}(\theta) \right., \quad j = 1, \dots, h.$$

Hence,  $\mathbb{G}_n(A) \equiv \mathbb{G}_n(A, \hat{\theta}_n) = \frac{1}{n} \sum_{j=1}^h 1_A(t_j) \hat{p}_j$  is Vardi's NPMLE evaluated at  $\hat{\theta}_n$ , which exists uniquely if  $\mathbf{G}(\hat{\theta}_n)$  is strongly connected. The semiparametric MLE  $\mathbb{G}_n$  is written explicitly as

$$\mathbb{G}_n(A) = \frac{\frac{1}{n} \sum_{j=1}^h 1_A(t_j) r_j \left[ \sum_{k=1}^s \lambda_{nk} w_{kj}(\hat{\theta}_n) \mathbb{V}_{nk}^{-1}(\hat{\theta}_n) \right]^{-1}}{\frac{1}{n} \sum_{j=1}^h r_j \left[ \sum_{k=1}^s \lambda_{nk} w_{kj}(\hat{\theta}_n) \mathbb{V}_{nk}^{-1}(\hat{\theta}_n) \right]^{-1}}.$$

The semiparametric MLE can also be written in terms of the empirical measure  $\mathbb{F}_n$ , defined by

$$(1.14) \quad \mathbb{F}_n(A) = \frac{1}{n} \sum_{i=1}^s \sum_{j=1}^{n_i} 1_A(Y_{ij}) = \sum_{i=1}^s \lambda_{ni} \frac{1}{n_i} \sum_{j=1}^{n_i} 1_A(Y_{ij}) = \sum_{i=1}^s \lambda_{ni} \mathbb{F}_{ni}(A)$$

for  $A \in \mathbf{B}$ ,  $i = 1, \dots, s$ , where  $Y_{ij}$  is the  $j$ th observation from sample  $i$ . We have

$$(1.15) \quad \mathbb{G}_n(A) = \frac{\int_A \left[ \sum_{k=1}^s \lambda_{nk} w_k(y, \hat{\theta}_n) \mathbb{V}_{nk}^{-1}(\hat{\theta}_n) \right]^{-1} d\mathbb{F}_n(y)}{\int_{\mathbf{Y}} \left[ \sum_{k=1}^s \lambda_{nk} w_k(y, \hat{\theta}_n) \mathbb{V}_{nk}^{-1}(\hat{\theta}_n) \right]^{-1} d\mathbb{F}_n(y)}.$$

The above procedure 1–3 is computationally attractive because it only requires calculation of Vardi's NPMLE  $\mathbb{G}_n$  one time; once  $(\hat{\theta}_n, \mathbb{V}_n)$  is obtained,  $\hat{p}$  is obtained through substitution only. Thus, in essence, the procedure only requires maximizing a function over a finite-dimensional parameter space and solving a system of equations of fixed  $(s - 1)$  dimension.

This paper is organized as follows. For identifiable  $s$ -sample semiparametric biased sampling models,  $s \geq 2$ , with  $w_s$  a constant function, large sample properties of the MLE are developed in Sections 2 through 6. We begin in Section 2 by discussing conditions for unique estimability of the model in the limit with probability 1. Section 3 provides heuristic discussion of the approach taken and defines notation. In Section 4, information bounds for estimation of  $\theta$  and  $G$  are calculated. Then results parallel to GVW's results on consistency, asymptotic normality and efficiency are given in Section 5. An asymptotically consistent estimator of the limiting covariance process is constructed by substitution of the MLE into the inverse generalized Fisher information. In Section 6, the theorems derived in this paper are applied to the examples introduced in Section 1. Proofs are presented in Section 7. Concluding comments and open problems are discussed in Section 8.

**2. Estimability and uniqueness.** Under what conditions does the likelihood (1.10) have a unique maximum in  $(\theta, G)$  with probability 1 as  $n \rightarrow \infty$ ? As stated in Theorem 1.2, the problem of maximizing the likelihood is equivalent to maximizing the partial likelihood (1.11), which can be accomplished by maximizing the profile partial likelihood (1.13). Thus our approach is to identify mild conditions under which the profile partial likelihood (1.13), and hence the full likelihood (1.10), has a unique maximum with probability 1.

For known  $\theta$ , Vardi [(1985), Lemma, page 197] showed that if the graph  $\mathbf{G}^*(\theta)$  is connected, then a unique solution  $\mathbb{V}_n(\theta)$  of (1.12) exists with probability 1 as  $n \rightarrow \infty$ . Therefore, assuming the graph  $\mathbf{G}^*(\theta)$  is connected for all  $\theta$  in some neighborhood  $\Theta_0 \subset \Theta$  of  $\theta_0$ , the first step in maximizing the profile partial likelihood always works when the sample size is large. Then the profile partial likelihood will have a unique maximum in  $\theta$  with probability 1 as  $n$  grows large if the limiting log profile partial likelihood is strictly concave on  $\Theta$ . Thus we have proved the following theorem.



**THEOREM 2.1** (Unique estimability via maximum partial likelihood estimation). *Suppose  $s \geq 2$ ,  $w_s$  is a constant function, the identifiability conditions of Theorem 1.1 hold, and for each  $i = 1, \dots, s$ ,  $n_i \rightarrow \infty$  and  $\lambda_{ni} \rightarrow \lambda_i > 0$ . Further suppose the limiting log profile partial likelihood is strictly concave on  $\Theta$ . Then with probability converging to 1, the likelihood (1.10) has a unique maximum, which can be obtained by the maximum partial likelihood estimation procedure 1–3.*

Henceforth assume the data and weights satisfy the hypotheses of Theorem 2.1, so that the  $s$ -sample biased sampling model is identifiable and uniquely estimable asymptotically. For this class of biased sampling models, we establish the large sample theory.

**3. Preliminaries: asymptotic theory for the MLE.** GVW characterized minimal assumptions on the known selection weight functions and the data under which Vardi's NPML  $\mathbb{G}_n$  is uniformly  $\sqrt{n}$ -consistent, asymptotically normal and efficient, where the uniformity is over a Donsker class. Essentially, the hypotheses needed for these results are connectedness of the graph  $\mathbf{G}^*$ ,  $n_i \rightarrow \infty$  with  $\lambda_{ni} \rightarrow \lambda_i > 0$ ,  $i = 1, \dots, s$ , the normalizing constants  $W_i(G)$  are finite and at least one weight function is bounded away from zero. Under these conditions (assuming that  $\mathbf{G}^*(\theta)$  is connected for all  $\theta$  in some neighborhood  $\Theta_o \subset \Theta$  of  $\theta_0$ ), we give sufficient additional integrability and smoothness hypotheses on the weight functions in order that the joint MLE  $(\hat{\theta}_n, \mathbb{G}_n)$  retains these asymptotic properties. The basic extra conditions are twice-differentiability of each weight function  $w_i(\cdot, \theta)$  in  $\theta$  at  $\theta_0$  and square-integrability of each "score for  $w_i$ "  $\partial/\partial\theta \log\{w_i(\cdot, \theta)\}$  in a neighborhood of  $\theta_0$ .

Consistency is derived by following GVW's approach and utilizing special features of the maximum profile partial likelihood estimation procedure. At first thought, since the dimension of the parameter space of the semiparametric estimate based on (1.10) grows with the sample size, one might expect that the estimate of  $\theta$  suffers from inconsistency due to infinitely many incidental parameters; see, for example, Kiefer and Wolfowitz (1956). However, the equivalence of maximizing the full likelihood over an infinite-dimensional parameter space with maximizing the partial likelihood over a finite-dimensional parameter space (Theorem 1.2) implies that this is not the case. Asymptotic normality of  $(\hat{\theta}_n, \mathbb{G}_n)$  is proved as follows. First, asymptotic normality of  $\hat{\theta}_n$  is established through application of a Master theorem [see, e.g., Bickel, Klaassen, Ritov and Wellner (1993), Theorem 1, page 312] to the finite collection of estimating equations formed by setting derivatives of the log profile partial likelihood (1.13.) to zero. The key stochastic, approximation hypothesis of the Master theorem is verified using tools in empirical process theory. Second, follow Gill, Vardi and Wellner's (1988) (henceforth GVW) proof of asymptotic normality of  $\mathbb{G}_n$ , adding the necessary hypotheses for it to go through for  $\mathbb{G}_n$  evaluated at  $\hat{\theta}_n$ . Efficiency is established by comparison of the limiting covariance structure of the MLE to the information bounds. The needed semi-

parametric efficiency theory is largely taken from Bickel, Klaassen, Ritov and Wellner (1993), henceforth referred to as BKRW.

An alternative method for establishing asymptotic normality and efficiency of infinite-dimensional maximum likelihood estimators, developed by Van der Vaart (1995), does not require calculation of the information bounds. However, the hypotheses of Van der Vaart’s powerful approach are difficult to verify for this problem. We have successfully applied Van der Vaart’s method for the case  $\mathbf{Y} = R$ , but have not yet succeeded for a general space  $\mathbf{Y}$  [Gilbert (1996)]. The difficulty is in identifying a set of functions on which the score operator acts and a metric space for this set for which Van der Vaart’s Fréchet differentiability and continuous invertibility hypotheses can both be verified. To obtain maximum generality with minimal assumptions, we directly extend the approach of GVW, exploiting the elegant structure in the maximum profile partial likelihood estimation procedure.

3.1. *Notation.* We give some notation, which corresponds to that used by GVW. Let  $(\theta_0, \underline{V}_0, G)$  denote the “true” parameter value, and let  $F_{i\theta}(A) \equiv F_i(A, \theta, G)$ ,  $A \in \mathbf{B}$ . Define

$$\tilde{w}_i(y, \theta) \equiv w_i(y, \theta)/W_i(\theta, G) \quad \text{and} \quad r(y, \theta) \equiv \left[ \sum_{i=1}^s \lambda_i \tilde{w}_i(y, \theta) \right]^{-1}.$$

Let  $\underline{r} = (r_1, \dots, r_s)^T$ , with  $r_i(y, \theta) \equiv \lambda_i \tilde{w}_i(y, \theta) / \sum_{k=1}^s \lambda_k \tilde{w}_k(y, \theta) = \lambda_i \tilde{w}_i(y, \theta) r(y, \theta)$ . Notice that  $\sum_{i=1}^s r_i(y, \theta) = 1$  for all  $y$  and  $\theta$ . Put  $\phi(y, \theta) = \sum_{i=1}^s r_i(y, \theta) (\dot{w}_i(y, \theta)/w_i(y, \theta))$ , and let  $\underline{w} \equiv (w_1, \dots, w_s)^T$ ,  $\tilde{\underline{w}} \equiv (\tilde{w}_1, \dots, \tilde{w}_s)^T$ ,  $\underline{W}_0 \equiv (W_1(\theta_0, G), \dots, W_s(\theta_0, G))^T$ ,  $\underline{\mathbb{W}}_n \equiv (\mathbb{W}_{n1}, \dots, \mathbb{W}_{ns})^T$ ,  $\underline{V}(\theta) = (V_1(\theta), \dots, V_s(\theta))^T$  and  $\underline{\lambda}$  be the diagonal matrix with diagonal  $(\lambda_1, \dots, \lambda_s)^T$ . For  $\theta \in \Theta$ , define an  $s \times s$  matrix  $M(\theta)$  by

$$\begin{aligned} M(\theta) &= \underline{\lambda}^{-1} - \int r(y, \theta) \underline{\tilde{w}}(y, \theta) \underline{\tilde{w}}^T(y, \theta) dG(y) \\ (3.1) \quad &= \underline{\lambda}^{-1} - G(r(\theta) \underline{\tilde{w}}(\theta) \underline{\tilde{w}}^T(\theta)). \end{aligned}$$

As in GVW (page 1083), the matrix  $M(\theta)$  is singular and has rank  $s - 1$ . Let  $M^-(\theta)$  be a  $\{1, 2\}$ -generalized inverse of  $M(\theta)$ . The needed facts about  $M^-(\theta)$  are given in Lemmas 5.1 and 5.2 of GVW, which hold for each fixed  $\theta \in \Theta$ . Let  $r_n, r_{ni}, \phi_n, \underline{\lambda}_n$  and  $M_n(\theta)$  equal  $r, r_i, \phi, \underline{\lambda}$  and  $M(\theta)$ , respectively, with the  $\lambda_i$ ’s replaced with  $\lambda_{ni}$ ’s and the  $W_i$ ’s replaced with  $\mathbb{W}_{ni}$ ’s.

Let  $\overline{F}_\theta(y) \equiv \sum_{i=1}^s \lambda_i F_{i\theta}(y)$  be the average of the biased sampling distributions (1.3), and  $\overline{F}_{n\theta}(y) = \sum_{i=1}^s \lambda_{ni} F_{i\theta}(y)$ . Notice from (1.3) that  $\overline{F}_\theta$  and  $G$  are related by

$$(3.2) \quad d\overline{F}_\theta(y) = r^{-1}(y, \theta) dG(y) \quad \text{and} \quad dG(y) = r(y, \theta) d\overline{F}_\theta(y).$$

With the observations viewed as i.i.d. random variables  $X_k = (I_k, Y_k)$ ,  $k = 1, \dots, n$ ,  $\bar{F}_\theta(h(Y))$  is the expectation of  $h$  under the law of  $X_1$ . For a function  $h$  depending on  $x = (i, y)$ , define

$$\begin{aligned} P_{(\theta, G)}[h(X)] &= E_{(\theta, G)}[h(I, Y)] \\ &= \sum_{i=1}^s \lambda_i F_{i\theta}(h(i, Y)) = \sum_{i=1}^s \lambda_i G(h(i, Y) \tilde{w}_i(Y, \theta)). \end{aligned}$$

(Evidently,  $P_{(\theta, G)}[h(Y)] = \bar{F}_\theta[h(Y)]$ .) Similarly define the biased sampling empirical measure by  $\mathbb{P}_n[h(X)] = \sum_{i=1}^s \lambda_{ni} \mathbb{F}_{ni}[h(i, Y)]$ . Let  $P_0 = P_{(\theta_0, G)}$  be the “true” semiparametric biased sampling distribution. Let  $\|\cdot\|_1$  denote the  $L_1$  norm  $\|\underline{x}\|_1 = \sum_{i=1}^d |x_i|$  for  $\underline{x} \in R^d$ .

**4. Information bounds.** We compute expressions for the limiting covariance process of an efficient estimator of  $(\theta_0, G)$  in the semiparametric biased sampling model. The notion according to which it is efficient is that of a best regular estimator sequence; see, for example, BKRW [page 21]. An estimator sequence is considered to be asymptotically efficient relative to the tangent set  $\dot{\mathbf{P}}_0^0$  if, among the class of regular estimators, its limiting variance has the least dispersion. The limiting process with “least dispersion” is defined by the information bounds  $I^{-1}(\theta_0)$  for  $\theta_0$  and  $I_G^{-1}$  for  $G$ .

Let  $\dot{l}_\theta$  be the score for  $\theta$ , the derivative of the log biased sampling density  $p$  of (1.1) with respect to  $\theta$ , given by

$$\begin{aligned} \dot{l}_\theta(x) &= \frac{\dot{w}_i(y, \theta)}{w_i(y, \theta)} - E\left[\frac{\dot{w}_i(Y, \theta)}{w_i(Y, \theta)} \mid I = i\right] \\ (4.1) \quad &= \frac{\dot{w}_i(y, \theta)}{w_i(y, \theta)} - G\left(\frac{\dot{w}_i(Y, \theta)}{w_i(Y, \theta)} \tilde{w}_i(Y, \theta)\right), \end{aligned}$$

where  $\dot{w}_i(y, \theta) \equiv (\partial/\partial\theta')w_i(y, \theta')|_{\theta=\theta}$ . We refer to  $\dot{w}_i(\cdot, \theta)/w_i(\cdot, \theta) = (\partial/\partial\theta') \log\{w_i(\cdot, \theta')\}|_{\theta=\theta}$  as the *score for  $w_i$* , and denote the  $k$ th component of  $\dot{w}_i$ ,  $k = 1, \dots, d$ , by  $\dot{w}_{ik}(y, \theta) \equiv (\partial/\partial\theta'_k)w_i(y, \theta')|_{\theta=\theta}$ . For a tangent  $h \equiv (\partial/\partial\eta) \log g_\eta|_{\eta=0} \in L_2^0(G)$ , where  $L_2^0(G)$  is the space of functions with  $\int h^2 dG < \infty$  and  $\int h dG = 0$ , the *score operator for  $g$*  equals

$$\begin{aligned} \dot{l}_g h(x) &= \frac{\partial}{\partial\eta} \log dP_{(\theta, G_\eta)}|_{\eta=0} \\ (4.2) \quad &= h(y) - E[h(Y) \mid I = i] \\ &= h(y) - G(h(Y) \tilde{w}_i(Y, \theta)). \end{aligned}$$

Define a weighted average score for  $\theta$  over the  $s$  samples by  $\phi^*(y, \theta) = \sum_{i=1}^s r_i(y, \theta) \dot{l}_\theta(i, y)$ .

4.1. Information bound for  $\theta_0$ .

THEOREM 4.1 (Information bound for  $\theta_0$ ). *Suppose each weight function  $w_i(y, \theta)$  is differentiable in  $\theta$  at  $\theta_0$ . Further suppose:*

(I) *For some  $i \in \{1, \dots, s\}$ , there does not exist a vector  $\underline{b} \in R^d$ ,  $\underline{b} \neq 0$ , such that*

$$(4.3) \quad \underline{b}^T \left( \frac{\dot{w}_i(y, \theta_0)}{w_i(y, \theta_0)} - \left[ \phi^*(y, \theta_0) - G(\phi^*(\theta_0)) \right] \right. \\ \left. - G(\phi^*(\theta_0)) \tilde{w}^T(\theta_0) \right) (M^-(\theta_0))^T \left[ \tilde{w}(y, \theta_0) r(y, \theta_0) - G(\tilde{w}(\theta_0) r(\theta_0)) \right] \\ \equiv b^T \left( \frac{\dot{w}_i(y, \theta_0)}{w_i(y, \theta_0)} - \alpha_*(y, \theta_0) \right)$$

*is a.s. constant with respect to the law  $F_{i_0}$ . Then:*

A. *The unique efficient score function  $i_{\theta_0}^*$  for  $\theta_0$  is*

$$i_{\theta_0}^*(x) = \frac{\dot{w}_i(y, \theta_0)}{w_i(y, \theta_0)} - \alpha_*(y, \theta_0) - E \left[ \frac{\dot{w}_i(Y, \theta_0)}{w_i(Y, \theta_0)} - \alpha_*(Y, \theta_0) \mid I = i \right],$$

*where*

$$\alpha_*(y, \theta_0) \equiv \phi^*(y, \theta_0) - G(\phi^*(\theta_0)) \\ + G(\phi^*(\theta_0)) \tilde{w}^T(\theta_0) (M^-(\theta_0))^T \\ \times \left[ \tilde{w}(y, \theta_0) r(y, \theta_0) - G(\tilde{w}(\theta_0) r(\theta_0)) \right].$$

B. *The efficient information matrix  $I(\theta_0)$  for  $\theta_0$  is nonsingular, and is given by*

$$(4.4) \quad P_0 \left[ i_{\theta_0}^*(I, Y) i_{\theta_0}^{*T}(I, Y) \right] = P_0 \left\{ \frac{\dot{w}_I(Y, \theta_0)}{w_I(Y, \theta_0)} \frac{\dot{w}_I^T(Y, \theta_0)}{w_I(Y, \theta_0)} \right. \\ \left. - \phi(Y, \theta_0) \phi(Y, \theta_0)^T \right\} \\ - A(\theta_0) M^-(\theta_0) A(\theta_0)^T,$$

*where  $M^-(\theta_0)$  is a  $\{1, 2\}$ -generalized inverse of  $M(\theta_0)$  and*

$$(4.5) \quad A(\theta_0) = G \left\{ \left( \frac{\dot{w}_1(Y, \theta_0)}{w_1(Y, \theta_0)} - \phi(Y, \theta_0) \right) \tilde{w}_1(Y, \theta_0), \right. \\ \left. \dots, \left( \frac{\dot{w}_s(Y, \theta_0)}{w_s(Y, \theta_0)} - \phi(Y, \theta_0) \right) \tilde{w}_s(Y, \theta_0) \right\}.$$

Expression (4.3) and the efficient score and information are independent of the choice of  $\{1, 2\}$ -generalized inverse  $M^-(\theta_0)$ .

When  $s = 2$ , the inverse efficient information (covariance) matrix  $I^{-1}(\theta_0)$  equals

$$(4.6) \quad \frac{1}{\lambda_1 \lambda_2} \left[ G \left( \frac{\dot{w}_1(\theta_0)}{w_1(\theta_0)} \frac{\dot{w}_1^T(\theta_0)}{w_1(\theta_0)} \tilde{w}_1(\theta_0) r(\theta_0) \right) - \frac{G \left( \frac{\dot{w}_1(\theta_0)}{w_1(\theta_0)} \tilde{w}_1(\theta_0) r(\theta_0) \right) G \left( \frac{\dot{w}_1^T(\theta_0)}{w_1(\theta_0)} \tilde{w}_1(\theta_0) r(\theta_0) \right)}{G(\tilde{w}_1(\theta_0) r(\theta_0))} \right]^{-1}.$$

REMARK 4.1. The reason  $\tilde{l}_{\theta_0}^*$  and  $I(\theta_0)$  do not depend on the choice of  $\{1, 2\}$ -generalized inverse  $M^-(\theta_0)$  is because each row of  $G(\phi^*(\theta_0) \tilde{w}^T(\theta_0))$ , and the vector  $[r(y, \theta_0) \tilde{w}(y, \theta_0) - G(r(\theta_0) \tilde{w}(\theta_0))]$  is in the range of  $M$ , equal to  $\text{Range}(M) = \{x: \lambda^T x = 0\}$ . See Lemma 5.2 (iv) of GVW. In applications, a convenient choice of  $M^-(\theta_0)$  is

$$(4.7) \quad M^-(\theta_0) = \begin{pmatrix} M_{11}^{-1}(\theta_0) & 0 \\ 0 & 0 \end{pmatrix},$$

where  $M_{11}(\theta_0)$  is the upper-left  $s - 1 \times s - 1$  submatrix of  $M(\theta_0)$ . GVW [page 1080] proved that  $M_{11}(\theta_0)$  is nonsingular under connectivity of  $\mathbf{G}^*(\theta_0)$ .

REMARK 4.2. Note that by the Cauchy–Schwarz inequality, the information for  $\theta_0$  when  $s = 2$  and  $d = 1$  is positive if and only if  $\dot{w}_1(\cdot, \theta_0)/w_1(\cdot, \theta_0)$  is nondegenerate, which is true if and only if hypothesis (I) of Theorem 4.1 holds. This sheds light on the necessity of hypothesis (I).

REMARK 4.3. The information for estimation of the regression parameter  $\theta_0$  in the Cox proportional hazards model takes a similar form. Let  $Z$  denote a covariate and  $T$  a failure time, with no censoring. Under the proportional hazards model  $\lambda(t|z) = \lambda(t|0) \exp\{\theta_0^T z\}$ , the information for  $\theta_0$  is  $I(\theta_0) = E \text{Var}[Z|T]$ . This is derived in BKRW (pages 80–82) by computation of the efficient score via orthogonality calculations. This same approach is used in the proof of Theorem 4.1, and  $I(\theta_0)$  of (4.4) has the same structure as  $I(\theta_0)$  in the Cox model, equal to  $P_0 \text{Var}[\dot{w}_I(Y, \theta_0)/w_I(Y, \theta_0)] - \alpha_*(Y, \theta_0)|I]$ .

4.2. *Information bound for G.* Since  $G$  is an infinite-dimensional parameter, the information bound is an *inverse information covariance functional*  $I^{-1}(P_0|G, \mathbf{P}) \equiv I_G^{-1}: \mathbf{H} \times \mathbf{H} \rightarrow R$ . It is defined, as in BKRW [page 184], by  $I_G^{-1}(h_1, h_2) \equiv E[\tilde{l}_G(\pi_{h_1}) \tilde{l}_G(\pi_{h_2})]$ , where  $\tilde{l}_G(\pi_h)$  is the efficient influence function for estimation of  $G(h)$ . Here  $\pi_h$  is the projection map from  $l^\infty(\mathbf{H})$  to  $R$ , defined by  $\pi_h(G) = G(h) = \int h dG$ , where we now view probability distributions  $G$  as elements of  $l^\infty(\mathbf{H})$ , the Banach space of all bounded functions  $z: \mathbf{H} \rightarrow R$ , which is equipped with the supremum norm.

**THEOREM 4.2** (Information bound for  $G$ ). *Suppose the conditions of Theorem 4.1 hold. Additionally suppose  $r(y, \theta_0)$  and  $r^{-1}(y, \theta_0)$  are bounded in  $y$ . Then the inverse information covariance functional  $I_G^{-1}$  for estimation of  $G$  is*

$$\begin{aligned}
 I_G^{-1}(h_1, h_2) &= G([h_1 - G(h_1)][h_2 - G(h_2)]r(\theta_0)) \\
 &\quad + G([h_1 - G(h_1)]r(\theta_0)\tilde{w}^T(\theta_0))M^{-}(\theta_0) \\
 &\quad \times G([h_2 - G(h_2)]r(\theta_0)\tilde{w}(\theta_0)) \\
 (4.8) \quad &\quad + G(h_1\alpha_\star^T(\theta_0)) * I^{-1}(\theta_0) * G(h_2\alpha_\star(\theta_0)) \\
 &= I_{G^0}^{-1}(h_1, h_2) \\
 &\quad + G(h_1\alpha_\star^T(\theta_0)) * I^{-1}(\theta_0) * G(h_2\alpha_\star(\theta_0))
 \end{aligned}$$

For  $h_1, h_2 \in \mathbf{H}$ . The first three lines of (4.8) equal  $I_{G^0}^{-1}$ , the inverse information covariance functional for estimation of  $G$  when  $\theta_0$  is known. This functional was derived by GVW [pages 1089–1091] through computation of the inverse information operator for  $g$ ,  $(\dot{l}_g^T \dot{l}_g)^{-1}$ . The fourth line of (4.8) is the inflation of the covariance functional when  $\theta_0$  is unknown.

Since  $G([h - G(h)]r(\theta_0)\tilde{w}(\theta_0))$  and the rows of  $G(\phi(\theta_0)\tilde{w}^T(\theta_0))$  are in Range  $(M)$ , the information bound for  $G$  is independent of the choice of generalized inverse  $M^{-}(\theta_0)$ .

**REMARK 4.4.** If  $\theta_0$  is known and there is no biasing in the sampling,  $r$  and all the weights are unity, so that  $I_G^{-1}$  reduces to  $I_G^{-1}(h_1, h_2) = G(h_1 h_2) - G(h_1)G(h_2)$ . This is the information bound for estimation of a probability measure  $G$  from an independent, identically distributed sample from  $G$ . It is well known that the empirical measure is regular and has limiting covariance given by  $G(h_1 h_2) - G(h_1)G(h_2)$ , so that it is efficient. Therefore, the inverse information covariance functional  $I_G^{-1}$  matches this known special case, as it must.

**REMARK 4.5.** The assumption that  $r^{-1}(y, \theta_0)$  is bounded in  $y$ , which is equivalent to each of the weight functions  $w_i(y, \theta_0)$ ,  $i = 1, \dots, s$ , being bounded in  $y$ , is only used for calculation of the information bound for  $G$ . This technically unpleasant assumption is needed to ensure existence of the inverse information operator  $(\dot{l}_g^T \dot{l}_g)^{-1}$ .

We write down  $I_G^{-1}$  for two-sample semiparametric biased sampling models. We calculate  $M^{-}(\theta_0)$  of (4.7) to equal the  $2 \times 2$  matrix with  $\lambda_1[\lambda_2 G(\tilde{w}_1(\theta_0) r(\theta_0))]^{-1}$  as the upper-left element and zeros elsewhere. Then direct calculation gives

$$\alpha^\star(y, \theta) = r_1(y, \theta) \left[ \frac{\dot{w}_1(y, \theta)}{w_1(y, \theta)} - G\left(\frac{\dot{w}_1(\theta)}{w_1(\theta)} \tilde{w}_1(\theta) r(\theta)\right) \right] / G\left(\tilde{w}_1(\theta) r(\theta)\right)$$

and, with all functions evaluated at  $\theta_0$ ,

$$\begin{aligned}
 I_G^{-1}(h_1, h_2) &= G([h_1 - G(h_1)][h_2 - G(h_2)]r) \\
 &\quad + G([h_1 - G(h_1)]r\tilde{w}_1) \frac{\lambda_1}{\lambda_2 G(r\tilde{w}_1)} G([h_2 - G(h_2)]r\tilde{w}_1) \\
 &\quad + \left\{ \lambda_1 \left[ G\left(h_1 \frac{\dot{w}_1}{w_1} \tilde{w}_1 r\right) G(\tilde{w}_1 r) - G\left(\frac{\dot{w}_1}{w_1} \tilde{w}_1 r\right) G(h_1 \tilde{w}_1 r) \right] \right. \\
 (4.9) \quad &\quad \times \left[ G\left(h_2 \frac{\dot{w}_1^T}{w_1} \tilde{w}_1 r\right) G(\tilde{w}_1 r) - G\left(\frac{\dot{w}_1^T}{w_1} \tilde{w}_1 r\right) G(h_2 \tilde{w}_1 r) \right] \Big\} \\
 &\quad \times \left\{ \lambda_2 G(\tilde{w}_1 r) \left[ G\left(\frac{\dot{w}_1}{w_1} \frac{\dot{w}_1^T}{w_1} \tilde{w}_1 r\right) G(\tilde{w}_1 r) \right. \right. \\
 &\quad \left. \left. - G\left(\frac{\dot{w}_1}{w_1} \tilde{w}_1 r\right) G\left(\frac{\dot{w}_1^T}{w_1} \tilde{w}_1 r\right) \right] \right\}^{-1}.
 \end{aligned}$$

The consistent plug-in estimate of  $\text{Cov}(\mathbb{G}_n(h_1), \mathbb{G}_n(h_2))$  is given by the empirical version of (4.9).

To work out the asymptotic covariance between  $\sqrt{n}(\hat{\theta}_n - \theta_0)$  and  $\sqrt{n}(\hat{G}_n(h) - G(h))$  for an efficient estimator  $(\hat{\theta}_n, \hat{G}_n)$ , suppose  $\hat{\theta}_n$  and  $\hat{G}_n$  are each asymptotically linear, with efficient influence functions  $\tilde{l}_{\theta_0} = I^{-1}(\theta_0)l_{\theta_0}^*$  and  $\tilde{l}_G$ , respectively. Set  $Z_{\theta_0} \equiv \lim_n \sqrt{n}(\hat{\theta}_n - \theta_0)$  and  $Z(h) \equiv \lim_n \sqrt{n}(\hat{G}_n(h) - G(h))$ . As in BKRW (panel 5.5.27, page 216), since  $\tilde{l}_G(\pi_h) = \dot{l}_g(l_g^T \dot{l}_g)^{-1}(\pi_h) - l_{\theta_0}^{*T} I^{-1}(\theta_0) G([h - G(h)]a_*(\theta_0))$ , where  $l_{\theta_0}^*$  is orthogonal to  $\text{Range}(\dot{l}_g)$ ,

$$\begin{aligned}
 \text{Cov}(Z_{\theta_0}, Z(h)) &= P_0(\tilde{l}_{\theta_0}(X)\tilde{l}_G(\pi_h)(X)) \\
 (4.10) \quad &= P_0\left(\left[ I^{-1}(\theta_0)l_{\theta_0}^* \right] \left[ -l_{\theta_0}^{*T} I^{-1}(\theta_0) G([h - G(h)]a_*(\theta_0)) \right] \right) \\
 &= -I^{-1}(\theta_0) G([h - G(h)]a_*(\theta_0)) \\
 &= -I^{-1}(\theta_0) G(ha_*(\theta_0)).
 \end{aligned}$$

## 5. Asymptotic behavior of $(\hat{\theta}_n, \mathbb{G}_n)$ .

**5.1. Consistency.** Consistency of  $(\hat{\theta}_n, \mathbb{G}_n)$  is established in two steps. First, consistency of  $\hat{\theta}_n$  is proved under the hypotheses that the log limiting profile partial likelihood  $l_{1\text{pro}}(\theta)$  and the log profile partial likelihood  $l_{n1\text{pro}}(\theta|\underline{x}) \equiv \log L_{n1\text{pro}}(\theta|\underline{x})$  for all  $n$  large enough are strictly concave on  $\Theta$ . Since  $\mathbb{G}_n \equiv \mathbb{G}_n(\hat{\theta}_n)$  is Vardi's NPMLLE evaluated at  $\hat{\theta}_n$ , consistency of  $\mathbb{G}_n$  essentially follows by consistency of  $\hat{\theta}_n$  and GVWs proof that Vardi's NPMLLE is consistent.

From (1.13), the log profile partial likelihood  $l_{n1\text{ pro}}$  is given by

$$(5.1) \quad \begin{aligned} \frac{1}{n} l_{n1\text{ pro}}(\theta|\underline{x}) &= \sum_{i=1}^s \lambda_{ni} \frac{1}{n_i} \sum_{j=1}^{n_i} \log \left\{ \frac{w_{ij}(\theta) \mathbb{V}_{ni}^{-1}(\theta)}{\sum_{k=1}^s \lambda_{nk} w_{kj}(\theta) \mathbb{V}_{nk}^{-1}(\theta)} \right\} \\ &= \mathbb{P}_n \log \{ \lambda_n^{-1} r_{nI}(Y, \theta) \}. \end{aligned}$$

**THEOREM 5.1** (Consistency of  $\hat{\theta}_n$ ). *Suppose that:*

- (i) *The parameter space  $\Theta \subset R^d$  is an open convex set.*
- (ii) *Each weight function  $w_i$  is continuously differentiable in  $\theta$  in a neighborhood of  $\theta_0$ .*
- (iii) *For all  $n$  larger than some fixed  $N$ ,  $l_{n1\text{ pro}}(\theta|\underline{x})$  is strictly concave on  $\Theta$ .*

*Then as each  $n_i \rightarrow \infty$  with  $\lambda_{ni} \rightarrow \lambda_i > 0$ ,*

$$(5.2) \quad \begin{aligned} \frac{1}{n} l_{n1\text{ pro}}(\theta|\underline{x}) &\rightarrow_{\text{a.s.}} l_{1\text{ pro}}(\theta) \\ &= \sum_{i=1}^s \lambda_i F_{i\theta} \log \left\{ \frac{w_i(Y, \theta) V^{-1}(\theta)}{\sum_{k=1}^s \lambda_k w_k(Y, \theta) V_k^{-1}(\theta)} \right\} \\ &= P_{(\theta, G)} \log \{ \lambda_I^{-1} r_I(Y, \theta) \} \end{aligned}$$

*for each fixed  $\theta \in \Theta$ . By strict concavity of  $l_{n1\text{ pro}}(\theta|\underline{x})$  on  $\Theta$  for all  $n > N$  and strict concavity of  $l_{1\text{ pro}}(\theta)$  on  $\Theta$ , which is implied by (iii) and (5.2), it follows that  $l_{n1\text{ pro}}(\theta|\underline{x})$  has a unique maximum at  $\hat{\theta}_n$  for every  $n > N$  and  $l_{1\text{ pro}}$  has a unique maximum at  $\theta_0$ , so that  $\hat{\theta}_n \rightarrow_{\text{a.s.}} \theta_0$  as  $n \rightarrow \infty$ .*

**REMARK 5.1.** Through theoretical results and simulations, Gilbert (1996) and Gilbert, Lele and Vardi (1999) (see especially Figure 2 and Theorems 5 and 6) showed that  $l_{n1\text{ pro}}(\theta|\underline{x})$  and  $l_{1\text{ pro}}(\theta)$  are strictly concave on  $\Theta \subset R^d$  for a large class of weight functions and data sets. A sufficient general condition for the concavity hypothesis (iii) to hold is that condition (I) of Theorem 4.1 holds for the function of (4.3) evaluated at every  $\theta \in \Theta$ . In the proof of Theorem 5.3, we show that the second derivative matrix of  $l_{1\text{ pro}}(\theta)$  equals  $-I(\theta)$ , so that positive definiteness of  $I(\theta)$  for all  $\theta \in \Theta$  implied by (I) for all  $\theta \in \Theta$  yields strict concavity of  $l_{1\text{ pro}}(\theta)$  on  $\Theta$ . Now, the second derivative matrix of  $l_{n1\text{ pro}}(\theta|\underline{x})$  is asymptotically equivalent to minus the observed information matrix  $-I_n(\theta)$  by regularity, and  $I_n(\theta)$  converges in probability to  $I(\theta)$  under the hypotheses of Proposition 5.1. Thus, under these hypotheses and condition (I) for all  $\theta \in \Theta$ , hypothesis (iii) holds. Condition (I) for all  $\theta$  in  $\Theta$  will usually hold if the following easy-to-verify condition holds: for some  $i \in \{1, \dots, s\}$  and all  $l \in \{1, \dots, d\}$  and every  $\theta \in \Theta$ ,  $\dot{w}_{il}(y, \theta)/w_i(y, \theta)$  and  $w_i(y, \theta)$  are not a.s. constant with respect to the law  $F_{i\theta}$ .



The first task in establishing consistency of  $\mathbb{G}_n$  is establishing consistency of  $\underline{\mathbb{V}}_n(\hat{\theta}_n)$ , which is defined as the solution of (1.12) evaluated at  $\hat{\theta}_n$ , or equivalently as the maximizer jointly with  $\hat{\theta}_n$  of the partial likelihood (1.11).

**PROPOSITION 5.1** [Consistency of  $\underline{\mathbb{V}}_n(\hat{\theta}_n)$  and  $\underline{\mathbb{W}}_n(\hat{\theta}_n)$ ]. *Suppose conditions (i)–(ii) of Theorem 5.1 hold, and its conclusion  $\hat{\theta}_n \xrightarrow{\text{a.s.}} \theta_0$ . In addition, suppose there exists a  $\delta > 0$  such that*

$$(iv) \sup_{\|\theta - \theta_0\|_1 < \delta} G(|\dot{w}_{ik}(\theta)/w_i(\theta)|\bar{w}_i(\theta)) < \infty \text{ for all } i = 1, \dots, s, k = 1, \dots, d.$$

*Then the equations (1.12) have, with probability 1 as  $n \rightarrow \infty$  with  $\lambda_{ni} \rightarrow \lambda_i > 0$ , the unique solution  $\underline{\mathbb{V}}_n(\hat{\theta}_n) = (\mathbb{V}_{n1}(\hat{\theta}_n), \dots, \mathbb{V}_{ns-1}(\hat{\theta}_n), 1)^T$  that satisfies*

$$\underline{\mathbb{V}}_n(\hat{\theta}_n) \xrightarrow{\text{a.s.}} \underline{V}(\theta_0) \equiv \underline{V}_0.$$

*Moreover,  $\underline{\mathbb{W}}_n(\hat{\theta}_n) \xrightarrow{\text{a.s.}} \underline{W}(\theta_0)$ .*

As in GVW [page 1081], let  $h_e$  be a fixed nonnegative  $G$ -integrable function, let  $\mathbf{C}$  be a Vapnik–Chervonenkis class of subsets of the probability space  $\mathbf{Y}$ , and consider the collection of functions

$$(5.3) \quad \mathbf{H} \equiv \{h_e 1_C : C \in \mathbf{C}\}.$$

We give conditions under which  $\mathbb{G}_n$  is a consistent estimator of  $G$  uniformly over  $\mathbf{H}$ .

**THEOREM 5.2** (Consistency of  $\mathbb{G}_n$ ). *Suppose the conditions of Proposition 5.1 hold, and that  $\mathbf{H}$  is defined as in (5.3), with  $h_e \in L_1(G)$ . Then as  $n \rightarrow \infty$  with  $\lambda_{ni} \rightarrow \lambda_i > 0$ ,*

$$(5.4) \quad \|\mathbb{G}_n - G\|_{\mathbf{H}} \equiv \sup_{h \in \mathbf{H}} \{|\mathbb{G}_n(h) - G(h)|\} \xrightarrow{\text{a.s.}} 0.$$

Theorem 5.2 has two immediate corollaries.

**COROLLARY 5.1.** *Assume that  $\mathbf{Y} = R^k$  and  $\mathbf{C}$  is a Vapnik–Chervonenkis class of subsets of  $\mathbf{Y}$ , and the conditions of Proposition 5.1 hold. Then as  $n \rightarrow \infty$  with  $\lambda_{ni} \rightarrow \lambda_i > 0$ ,*

$$(5.5) \quad \|\mathbb{G}_n - G\|_{\mathbf{C}} \equiv \sup_{C \in \mathbf{C}} |\mathbb{G}_n(C) - G(C)| \xrightarrow{\text{a.s.}} 0.$$

**COROLLARY 5.2.** *Suppose  $G(|h|) < \infty$ , and the conditions of Proposition 5.1 hold. Then  $\mathbb{G}_n(h) \xrightarrow{\text{a.s.}} G(h)$  as  $n \rightarrow \infty$  with  $\lambda_{ni} \rightarrow \lambda_i > 0$ .*

5.2. *Asymptotic distributions.* We formulate the asymptotic distribution theory for  $\sqrt{n}(\mathbb{G}_n - G)$  in terms of the  $s$ -sample  $Y$ -marginal empirical process  $\mathbb{X}_n$ , defined by

$$\mathbb{X}_n \equiv \sqrt{n}(\mathbb{F}_n - \overline{F}_{n0}) = \sum_{i=1}^s \sqrt{\lambda_{ni}} \sqrt{n_i} (\mathbb{F}_{ni} - F_{i0}),$$

where  $\mathbb{F}_n$  and  $\mathbb{F}_{ni}$  are as in (1.14). If  $\mathbf{F}$  is a Donsker class for each  $F_{i0}$ ,  $i = 1, \dots, s$ , then  $\{\mathbb{X}_n(f): f \in \mathbf{F}\}$  converges in distribution in  $l^\infty(\mathbf{F})$  to the mean zero Gaussian process  $\{\mathbb{X}(f): f \in \mathbf{F}\}$  with covariance

$$\begin{aligned} \text{Cov}(\mathbb{X}(h_1), \mathbb{X}(h_2)) &= \sum_{i=1}^s \lambda_i \{F_{i0}(h_1 h_2) - F_{i0}(h_1) F_{i0}(h_2)\} \\ &= G(r^{-1}(\theta_0) h_1 h_2) - G(h_1 \underline{w}^T(\theta_0)) \underline{\lambda} G(h_2 \underline{w}(\theta_0)). \end{aligned}$$

To avoid problems with measurability, convergence in distribution is defined using outer expectations as in Dudley (1985).

For fixed  $\theta \in \Theta$ , define the *biased sampling empirical process* by  $\mathbb{Z}_n(\cdot, \theta) \equiv \sqrt{n}(\mathbb{G}_n(\theta) - G)$ , regarded as a process indexed by a collection of functions  $\mathbf{H} \subset L_2(G)$ . Thus, for  $h \in \mathbf{H}$ ,

$$\mathbb{Z}_n(h, \theta) = \int h d\mathbb{Z}_n(\theta) = \sqrt{n} \int h d(\mathbb{G}_n(\theta) - G).$$

We formulate a limit theorem for  $\mathbb{Z}_n(\cdot, \hat{\theta}_n)$  jointly with the  $\sqrt{n}(\hat{\theta}_n - \theta_0)$  process. This joint process, indexed by  $\Theta \times \mathbf{H}$ , is defined by

$$\mathbb{Z}'_n(h, \hat{\theta}_n) = \left( \sqrt{n}(\hat{\theta}_n - \theta_0), \mathbb{Z}_n(h, \hat{\theta}_n) \right)^T.$$

The appropriate limiting Gaussian distribution is given by

$$(5.6) \quad \mathbb{Z}'(h) = \left( Z_{\theta_0}, -G(h \alpha_\star^T(\theta_0)) * Z_{\theta_0} + Z(h) \right)^T,$$

where  $Z_{\theta_0} \equiv \lim_n \sqrt{n}(\hat{\theta}_n - \theta_0)$  and  $Z$  is a mean zero Gaussian process defined in (2.18) of GVW, with covariance function equal to  $I_{G^0}^{-1}$ , displayed as the first three lines of expression (4.8). The process  $\mathbb{Z}'$  has covariance function

$$(5.7) \quad \begin{aligned} &\text{Cov}(\mathbb{Z}'(h_1), \mathbb{Z}'(h_2)) \\ &= \begin{bmatrix} I^{-1}(\theta_0) & -I^{-1}(\theta_0) * G(h_2 \alpha_\star(\theta_0)) \\ -G(h_1 \alpha_\star^T(\theta_0)) * I^{-1}(\theta_0) & I_G^{-1}(h_1, h_2) \end{bmatrix} \end{aligned}$$

for  $h_1, h_2 \in \mathbf{H}$ , where  $\alpha_\star$  and  $I^{-1}(\theta_0)$  are as in Theorem 4.1 and  $I_G^{-1}$  is as in Theorem 4.2.

As in GVW, take  $\mathbf{H}$  to be a collection of functions with  $G$ -integrable envelope function  $h_e$  such that  $G(h_e^2 r(\theta_0)) < \infty$  and the class of functions  $\mathbf{F} =$

$\{hr(\theta_0): h \in \mathbf{H}\}$  is Donsker for each  $F_{i0}$ . We also require that the weight functions (in a neighborhood of  $\theta_0$ ) be included in  $\mathbf{H}$ , that is, for some  $\delta > 0$ ,

$$(5.8) \quad \mathbf{W}_\delta \equiv \bigcup_{\|\theta - \theta_0\|_1 < \delta} \{w_i(\theta)r(\theta), \dots, w_s(\theta)r(\theta)\}$$

is Donsker for each  $F_{i0}$ . Additionally, assume that the finite collection of score functions

$$(5.9) \quad \mathbf{W} \equiv \left\{ \frac{\dot{w}_{1k}(y, \theta_0)}{w_{1k}(y, \theta_0)}, \dots, \frac{\dot{w}_{sk}(y, \theta_0)}{w_{sk}(y, \theta_0)} : k = 1, \dots, d \right\}$$

is Donsker for each  $F_{i0}$ .

Before stating our limit theorem for  $\mathbb{Z}'_n$ , we establish asymptotic normality (and efficiency) of  $\hat{\theta}_n$ . We apply the Master theorem, originating from Huber (1967), which applies to estimators which approximately zero out a collection of score or estimating equations. In our case, the estimating equations arise from setting the derivative of the log profile partial likelihood (5.1) equal to zero. Recall that  $n^{-1}l_{n1\text{pro}}(\theta|\underline{x}) = \mathbb{P}_n[\log\{\lambda_n^{-1}r_{nI}(Y, \theta, \underline{V}_n(\theta))\}]$ . Let  $\psi_\theta(x) \equiv \psi(x, \theta, \underline{V}(\theta))$  equal  $(\partial/\partial\theta') \log\{\lambda_i^{-1}r_i(y, \theta', \underline{V}(\theta'))\}|_\theta$ . By the chain rule, we calculate

$$(5.10) \quad \begin{aligned} \psi_\theta(x) &= \frac{\partial}{\partial\theta'} \log\{\lambda_i^{-1}r_i(y, \theta', \underline{V})\} \Big|_{(\theta, \underline{V}(\theta))} \\ &+ \frac{\partial}{\partial\theta'} \underline{V}^T(\theta') \Big|_\theta \frac{\partial}{\partial\underline{V}} \log\{\lambda_i^{-1}r_i(y, \theta, \underline{V})\} \Big|_{(\theta, \underline{V}(\theta))} \\ &= \left[ \frac{\dot{w}_i(y, \theta)}{w_i(y, \theta)} - \phi(y, \theta) \right] + \frac{\partial}{\partial\theta'} \underline{V}^T(\theta') \Big|_\theta * \underline{V}^{-1}(\theta)[-e_i + \underline{r}(y, \theta)], \end{aligned}$$

where  $e_i$  is the  $s$ -column vector with a 1 at position  $i$  and zeros elsewhere.

The matrix  $(\partial/\partial\theta') \underline{V}^T(\theta')|_\theta$  can be calculated by differentiating both sides of the profile likelihood equations defining  $\underline{V}(\theta)$ ,  $P[w_i(\theta)V_i^{-1}(\theta)/\sum_{k=1}^s \lambda_k w_k(\theta)V_k^{-1}(\theta)] = 1, i = 1, \dots, s-1$ . Setting  $\underline{d}_i(\theta) = (\partial/\partial\theta')V_i(\theta')|_\theta/V_i(\theta), i = 1, \dots, s$  (with  $\underline{d}_s(\theta) = 0$ ), this gives

$$\begin{aligned} \underline{d}_i(\theta) &= G \left[ \left( \frac{\dot{w}_i(\theta)}{w_i(\theta)} - \phi(\theta) \right) \tilde{w}_i(\theta) \right] \\ &+ \sum_{k=1}^s \lambda_k \underline{d}_k(\theta) G(r(\theta)\tilde{w}_i(\theta)\tilde{w}_k(\theta)), \quad i = 1, \dots, s. \end{aligned}$$

In matrix notation this is written, with  $\underline{d}(\theta) = [\underline{d}_1(\theta), \dots, \underline{d}_s(\theta)]$ , as

$$\underline{d}(\theta) = A(\theta) + \underline{d}(\theta)\underline{\lambda}G(r(\theta)\underline{\tilde{w}}(\theta)\underline{\tilde{w}}^T(\theta)).$$

By definition of the matrix  $M$ , we find that the solution  $\underline{d}(\theta)$  is given by  $A(\theta)M^-(\theta)\underline{\lambda}^{-1}$ . Substituting this into (5.10) gives

$$(5.11) \quad \psi_\theta(x) = \frac{\dot{w}_i(y, \theta)}{w_i(y, \theta)} - \phi(y, \theta) + A(\theta)M^-(\theta)\underline{\lambda}^{-1}[-e_i + r(y, \theta)].$$

Note that, if the profile partial likelihood estimator  $\hat{\theta}_n$  is efficient, the function  $\psi_\theta$  will equal the efficient score function  $\dot{l}_\theta^*$  defined in Theorem 4.1.

Letting  $\psi_{n\theta}(x) = \psi_n(x, \theta, \underline{V}_n(\theta))$  be the empirical version of  $\psi_\theta(x)$ , with the  $\lambda_i$ ,  $r_i$  and  $V_i$  replaced with  $\lambda_{ni}$ ,  $r_{ni}$  and  $V_{ni}$ , respectively, it follows that

$$(5.12) \quad \begin{aligned} \Psi_n(\theta|\underline{x}) &\equiv \frac{\partial}{\partial \theta'} \frac{1}{n} l_{n1 \text{ pro}}(\theta'|\underline{x})|_\theta \\ &= \mathbb{P}_n \left[ \frac{\partial}{\partial \theta'} \left[ \lambda_{n1}^{-1} r_{n1}(Y, \theta', \underline{V}_n(\theta')) \right] \Big|_\theta \right] \\ &= \mathbb{P}_n \left[ \psi_n(X, \theta, \underline{V}_n(\theta)) \right] \\ &= \mathbb{P}_n [\psi_\theta(X)] + o_p(n^{-1/2}), \end{aligned}$$

where the last equality will be verified in the proof of Theorem 5.3.

By the law of large numbers,  $\mathbb{P}_n \psi_\theta(X) \rightarrow_p P_{(\theta, G)} \psi_\theta(X)$ . Let  $\dot{\Psi}_0 = P_0 \nabla_\theta \psi_\theta(X)|_{\theta_0}$ , the derivative matrix of  $P_0 \psi_\theta(X)$  evaluated at  $\theta_0$ .

**THEOREM 5.3** (Asymptotic normality of  $\hat{\theta}_n$ ). *Suppose there exists a  $\delta > 0$  such that  $\mathbf{W}_\delta$  and  $\mathbf{W}$  of (5.8) and (5.9) are Donsker for each  $F_{i0}$ . Further assume hypothesis (I) of Theorem 4.1 and the hypotheses of Proposition 5.1. Additionally suppose:*

- (v) *Each weight function is twice-differentiable in  $\theta$  at  $\theta_0$ .*
- (vi) *There exists a  $\delta > 0$  such that  $\sup_{\|\theta - \theta_0\|_1 < \delta} r(y, \theta)$  is bounded in  $y$ .*
- (vii) *There exists a  $\delta > 0$  such that  $\sup_{\|\theta - \theta_0\|_1 < \delta} G(|\dot{w}_{ik}(\theta)/w_i(\theta)|^2 \tilde{w}_i(\theta_0)) < \infty$  for  $i = 1, \dots, s, k = 1, \dots, d$ .*
- (viii) *For  $i = 1, \dots, s - 1, k = 1, \dots, d$ , there exists a constant  $K$ , a  $\delta > 0$  and an  $\alpha > d/2$  such that*

$$\sup_{y \in \mathbf{Y}} \left[ \left| \frac{\dot{w}_{ik}(y, \theta_1)}{w_i(y, \theta_1)} - \frac{\dot{w}_{ik}(y, \theta_2)}{w_i(y, \theta_2)} \right| \right] \leq K \|\theta_1 - \theta_2\|_1^\alpha$$

for all  $\theta_1, \theta_2 \in \Theta_\delta \equiv \{\theta \in \Theta: \|\theta - \theta_0\|_1 < \delta\}$ .

Then  $\dot{\Psi}_0$  is nonsingular and  $\hat{\theta}_n$ , the unique maximizer of the profile likelihood (1.13), satisfies

$$\begin{aligned} \sqrt{n}(\hat{\theta}_n - \theta_0) &= -\frac{1}{\sqrt{n}} \sum_{i=1}^n \dot{\Psi}_0^{-1} \psi_0(X_i) + o_p(1) \\ &\rightarrow_d N_d(\underline{0}, \Sigma), \end{aligned}$$

where  $\Sigma = \Psi_0^{-1} P_0[\psi_0(X)\psi_0^T(X)](\Psi_0^{-1})^T$ . Here  $\dot{\Psi}_0 = -P_0[\psi_0(X)\dot{\psi}_0^T(X)]$ , so that  $\Sigma = -\dot{\Psi}_0^{-1}$ . Furthermore,  $-\dot{\Psi}_0 = I(\theta_0)$ , the efficient information matrix  $I(\theta_0)$  defined in (4.4), so that

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \rightarrow_d Z_{\theta_0} \sim N_d(\underline{0}, I^{-1}(\theta_0)).$$

REMARK 5.2. To establish asymptotic normality of  $\hat{\theta}_n$ , the present techniques of proof rely upon square  $F_{i0}$ -integrability of each score  $\dot{w}_i(\theta)/w_i(\theta)$  uniformly in  $\theta$  in a neighborhood of  $\theta_0$ . This condition is equivalent to  $\sup_{\|\theta-\theta_0\|_1 < \delta} P_0[\dot{l}_\theta^T(X)\dot{l}_\theta(X)] < \infty$ , where  $\dot{l}_\theta$  is the score for  $\theta$  defined in (4.1).

REMARK 5.3. Hypothesis (I) is used to establish invertibility of  $\dot{\Psi}_0$ . Empirical process theory is used to verify the stochastic, approximation condition of the Master theorem. It is implied if a certain class of functions is Donsker [see Van der Vaart (1995), Lemma 1], and is why the classes  $\mathbf{W}_\delta$  and  $\mathbf{W}$  are assumed to be Donsker. The uniform Lipschitz condition (viii) is used in verifying that classes of differences of scores for weights in a neighborhood of  $\theta_0$  are Donsker.

Given asymptotic normality of  $\hat{\theta}_n$ , asymptotic normality of the joint process  $Z'_n$  can be established.

THEOREM 5.4 (Asymptotic normality of  $(\hat{\theta}_n, \mathbb{G}_n)$ ). *Let  $\mathbf{H}$  be a collection of functions with envelope  $h_e$  such that  $\mathbf{F} = \{hr(\theta_0): h \in \mathbf{H}\}$  is Donsker for each  $F_{i0}$ ,  $G(h_e^2 r(\theta_0)) = \bar{F}_0(h_e^2 r^2(\theta_0)) < \infty$  and there exists  $\delta > 0$  such that  $\sup_{\|\theta-\theta_0\|_1 < \delta} G(h_e|\dot{w}_{ik}(\theta)/w_i(\theta)|^2 \tilde{w}_i(\theta_0)) < \infty$  for each  $i = 1, \dots, s, k = 1, \dots, d$ . Further assume the hypotheses of Theorem 5.3.*

*Then  $(\hat{\theta}_n, \mathbb{G}_n)$ , the solution of procedure (1)–(3), satisfies*

$$\begin{aligned} Z'_n(\cdot, \hat{\theta}_n) &= \left(\sqrt{n}(\hat{\theta}_n - \theta_0), Z_n(\cdot, \hat{\theta}_n)\right)^T \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n (\tilde{l}_{\theta_0}(X_i), \tilde{l}_G(X_i))^T + o_P(1) \\ &\Rightarrow Z' \text{ in } \Theta \times l^\infty(\mathbf{H}), \end{aligned}$$

where  $Z'$  is a tight mean zero Gaussian process in  $\Theta \times l^\infty(\mathbf{H})$  with covariance process given by (5.7).

Conditions for a class  $\mathbf{F}$  to be Donsker are well characterized; for example, see Pollard (1984) and Van der Vaart and Wellner (1996). A special case is given in the following corollary, with important application  $\mathbf{Y} = R^k, h_e = 1$ , and  $\mathbf{H}$  the collection of all indicator functions of lower-left orthants, or of all rectangles, or of all balls.

COROLLARY 5.3. *Suppose  $\mathbf{H} = \{h_e 1_C: C \in \mathbf{C}\}$ , where  $\mathbf{C}$  is a Vapnik–Chervonenkis class of sets. If  $G(h_e^2 r(\theta_0)) < \infty$  and there exists  $\delta > 0$  such*

that  $\sup_{\|\theta - \theta_0\|_1 < \delta} G(h_e | \dot{w}_{ik}(\theta) / w_i(\theta) |^2 \tilde{w}_i(\theta_0)) < \infty$ ,  $i = 1, \dots, s$ ,  $k = 1, \dots, d$ , then the result of Theorem 5.4 follows from the hypotheses of Theorem 5.3.

Notice that  $\mathbf{H}$  of the corollary equals  $\mathbf{H}$  of (5.3), the space over which the semiparametric MLE  $\mathbb{G}_n$  was proved to be uniformly consistent. For the distributional result, square  $G$ -integrability of the envelope  $h_e$  is required, while only  $G$ -integrability is needed for consistency.

Since the limiting covariance process coincides with the information bounds for  $\theta$  and  $G$ , with  $\text{AsymCov}(Z_{\theta_0}, Z(h))$  equal to covariance expression (4.10), the MLE  $(\hat{\theta}_n, \mathbb{G}_n)$  is asymptotically efficient relative to a tangent set  $\mathbf{P}_0^0$ . It is efficient in the following sense: for a given, fixed biased sampling problem as in (1.1)–(1.3) (the number of samples  $s \geq 2$ , the parametric form for the weight functions  $w_1, \dots, w_s$  and the sampling fractions  $\lambda_1, \dots, \lambda_s$  are all known), among the class of regular estimators, its limiting process has the least dispersion. See Remark 3.1 and Problem 6.1 in GVW for a discussion of designing the biased sampling in order to optimize some criteria.

REMARK 5.4. Note that  $G(w_i(\theta)r(\theta)) < \infty$  is automatically satisfied if  $W_i(\theta, G) < \infty$  and  $\lambda_i > 0$ , since  $G(w_i(\theta)r(\theta)) = [W_i(\theta, G)/\lambda_i]G(r_i(\theta)) \leq W_i(\theta, G)/\lambda_i$ . This, coupled with the hypothesis  $G(h_e^2 r(\theta_0)) < \infty$ , implies that the union  $\mathbf{F}_\delta \equiv \mathbf{F} \cup \mathbf{W}_\delta$  is Donsker because  $\mathbf{F}$  and  $\mathbf{W}_\delta$  are; see, for example, Van der Vaart and Wellner [(1996), Example 2.10.7]. This fact is used in the proof of Theorem 5.4.

REMARK 5.5 (Estimation of variability). The form of the covariance process (5.7) suggests a plug-in procedure for estimating it. Simply construct the sample analogue of (5.7), with  $\theta_0$ ,  $G$  and  $\lambda_i$  replaced by  $\hat{\theta}_n$ ,  $\mathbb{G}_n$  and  $\lambda_{ni}$  throughout. This estimator converges in probability and almost surely to (5.7) under natural conditions. This covariance estimator can also be obtained from the full likelihood (1.10) as follows. Form the Hessian of  $l_n = \log L_n$  by taking second derivatives with respect to the mass points of  $G$  (the  $p_j$ ) and  $\theta$ . Invert the Hessian, which is an  $(h + d) \times (h + d)$  matrix. For  $\underline{c} \in R^d$  and  $h \in \mathbf{H}$ , form the vector  $(c_1, \dots, c_d, h(u_1), \dots, h(u_h))$ , where  $u_j$ ,  $j = 1, \dots, h$ , is the location of the mass point  $p_j$ . Pre- and post-multiply the inverted Hessian to form the estimator of the asymptotic variance of  $\underline{c}^T \sqrt{n}(\hat{\theta}_n - \theta_0) + \mathbb{Z}_n(h, \hat{\theta}_n)$ . The latter method has been implemented in a computer program, and consistency corroborated in simulations [Gilbert, Lele and Vardi (1999)].

REMARK 5.6 (Bootstrap covariance estimator). We conjecture that the bootstrap covariance estimator is asymptotically first-order correct. We have verified this [Gilbert (1996)] for the case  $\mathbf{Y} = R$  by utilizing the fact that bootstrapped  $Z$ -estimators are first-order asymptotically correct if they satisfy a mild integrability condition and the hypotheses of Van der Vaart's (1995)  $Z$ -theorem [Wellner and Zhan (1998)]. Bootstrap confidence limits for  $(\hat{\theta}_n, \mathbb{G}_n)$  are studied in simulations by Gilbert, Lele and Vardi (1999).

REMARK 5.7 (Hypothesis tests of  $H_0: \theta_0 = 0$ ). Gilbert, Lele and Vardi (1999) developed and investigated finite-sample properties of Wald, efficient score and likelihood ratio tests of the hypothesis  $H_0: \theta_0 = 0$ . It is straightforward to verify that the statistics are asymptotically  $\chi_d^2$  under  $H_0$  if the hypotheses of Theorem 5.4 hold. In addition, the profile likelihood-based confidence set for  $\theta_0$  defined by  $\{\theta \in \Theta | 2[l_{n1\text{ pro}}(\hat{\theta}_n | \underline{x}) - l_{n1\text{ pro}}(\theta | \underline{x})] \leq \chi_{d, 1-\alpha}^2\}$  has correct converge probability  $1 - \alpha$  in the limit.

PROPOSITION 5.2 (Asymptotic normality of  $\underline{\mathbb{W}}_n$ ). *Suppose the hypotheses of Theorem 5.4 hold. Then*

$$\begin{aligned}
 \sqrt{n}(\underline{\mathbb{W}}_n - \underline{W}_0) &\rightarrow_d \mathbb{X}(r(\theta_0)(\underline{w}(\theta_0) - \underline{W}_0)) \\
 &+ G(r(\theta_0)(\underline{w}(\theta_0) - \underline{W}_0)\underline{\tilde{w}}^T(\theta_0))M^-(\theta_0)\mathbb{X}(r(\theta_0)\underline{\tilde{w}}(\theta_0)) \\
 (5.13) \quad &+ \underline{\underline{W}}_0 \left[ G(\underline{\tilde{w}}(\theta_0)\alpha_*^T(\theta_0)) \right. \\
 &\quad \left. + G\left(\frac{\dot{w}_1(\theta_0)}{w_1(\theta_0)}\tilde{w}_1(\theta_0), \dots, \frac{\dot{w}_s(\theta_0)}{w_s(\theta_0)}\tilde{w}_s(\theta_0)\right)^T \right] * Z_{\theta_0} \\
 &\sim N_s(\underline{0}, \Sigma),
 \end{aligned}$$

where

$$\begin{aligned}
 \Sigma &= G\left(r(\theta_0)(\underline{w}(\theta_0) - \underline{W}_0)(\underline{w}(\theta_0) - \underline{W}_0)^T\right) \\
 &+ G(r(\theta_0)(\underline{w}(\theta_0) - \underline{W}_0)\underline{\tilde{w}}^T(\theta_0))M^-(\theta_0) \\
 &\quad \times G(r(\theta_0)\underline{\tilde{w}}(\theta_0)(\underline{w}(\theta_0) - \underline{W}_0)^T) \\
 (5.14) \quad &+ \underline{\underline{W}}_0 \left[ G(\underline{\tilde{w}}(\theta_0)\alpha_*^T(\theta_0)) \right. \\
 &\quad \left. + G\left(\frac{\dot{w}_1(\theta_0)}{w_1(\theta_0)}\tilde{w}_1(\theta_0), \dots, \frac{\dot{w}_s(\theta_0)}{w_s(\theta_0)}\tilde{w}_s(\theta_0)\right)^T \right] * I^{-1}(\theta_0) \\
 &\quad \times \left[ G(\underline{\tilde{w}}(\theta_0)\alpha_*^T(\theta_0)) \right. \\
 &\quad \left. + G\left(\frac{\dot{w}_1(\theta_0)}{w_1(\theta_0)}\tilde{w}_1(\theta_0), \dots, \frac{\dot{w}_s(\theta_0)}{w_s(\theta_0)}\tilde{w}_s(\theta_0)\right)^T \right]^T \underline{\underline{W}}_0.
 \end{aligned}$$

Notice that when  $w_s$  is a constant function (say  $w_s = c$ ), the  $s$ th row and column of the covariance matrix  $\Sigma$  are filled with zeros. Thus, the above limiting expression captures the degeneracy in the  $s$ th weight function. Note that since  $\mathbb{V}_{ni} = \mathbb{W}_{ni}/\mathbb{W}_{ns} = c^{-1}\mathbb{W}_{ni}$ , it follows automatically that  $\sqrt{n}(\underline{\mathbb{V}}_n - \underline{V}_0) \rightarrow_d c^{-1}N_s(0, \Sigma)$ , where  $\Sigma$  is as in (5.14).

Finally, consider estimation of the biased distributions  $F_{10}, \dots, F_{s0}$ . The semiparametric MLE  $\hat{F}_{ni}$  of  $F_{i0}$  is given by

$$\hat{F}_{ni}(h, \hat{\theta}_n) = \mathbb{G}_n(hw_i(\hat{\theta}_n)) / \mathbb{G}_n(w_i(\hat{\theta}_n))$$

for  $h \in \mathbf{H}$ ,  $i = 1, \dots, s$ . Set  $\mathbb{Y}_{ni} = \sqrt{n}(\hat{F}_{ni}(\cdot, \hat{\theta}_n) - F_{i0})$ ,  $i = 1, \dots, s$ , and let  $\underline{\mathbb{Y}}_n = (\mathbb{Y}_{n1}, \dots, \mathbb{Y}_{ns})^T$ , which is a vector of processes indexed by a collection of functions  $\mathbf{H} \subset L_2(\bar{F}_0)$ . Define a vector of Gaussian processes  $\underline{Y}: \mathbf{H} \rightarrow R^s$  by

$$\begin{aligned} \underline{Y}(h) = & \mathbb{X}(hr(\theta_0)\underline{\tilde{w}}(\theta_0)) \\ (5.15) \quad & - \left[ G(h\underline{\tilde{w}}(\theta_0))\underline{\lambda}^{-1} - (hr(\theta_0)\underline{\tilde{w}}(\theta_0)\underline{\tilde{w}}^T(\theta_0)) \right] M^-(\theta_0)\mathbb{X}(r(\theta_0)\underline{\tilde{w}}(\theta_0)) \\ & + \left[ -G(h\underline{\tilde{w}}(\theta_0)a_\star^T(\theta_0)) + G\left(h\left[l_{\theta_0}(i, Y), \dots, l_{\theta_0}(s, Y)\right]^T\right) \right] * Z_{\theta_0} \end{aligned}$$

PROPOSITION 5.3 (Asymptotic normality of  $\underline{\mathbb{Y}}_n$ ). *Let  $\delta > 0$  and  $\mathbf{H}$  be a collection of functions such that  $\mathbf{F} \equiv \bigcup_{i=1}^s \{hr(\theta_0)w_i(\theta_0): h \in \mathbf{H}\}$  is Donsker and  $G(h_e^2 r(\theta)) < \infty$  for all  $\|\theta - \theta_0\|_1 < \delta$ . Suppose the hypotheses of Theorem 5.4 hold. Then  $\underline{\mathbb{Y}}_n \Rightarrow \underline{Y}$  in  $l^\infty(\mathbf{H})^s$  as  $n \rightarrow \infty$  with  $\lambda_{ni} \rightarrow \lambda_i > 0$ .*

### 6. Examples and applications.

EXAMPLE 6.1 (Examples 1.1 and 1.2 continued; generalized logistic regression model). For which GLR models is the MLE  $(\hat{\theta}_n, \mathbb{G}_n)$  consistent, asymptotically normal and efficient? For a large class of useful ones, as we now assert. First consider the univariate  $s$ -sample model (1.5) with  $\mathbf{Y} = R$ ,  $\mathbf{H} = \{1_{[0, t]}: t \in [0, \infty)\}$ . Take  $h_i(y, \theta) = \sum_{k=1}^d h_{ik}(y)\theta_k$  for known nonnegative functions  $h_{ik}$  with  $h_{ik}(0) = 0$ ,  $i = 1, \dots, s - 1$ , and set  $h_s \equiv 0$ . In this case,

$$\frac{\dot{w}_{ik}(y, \theta)}{w_i(y, \theta)} = h_{ik}(y)$$

and

$$w_i(y, \theta)r(y, \theta) = \frac{\exp\{\sum_{k=1}^d h_{ik}(y)\theta_k\}}{\sum_{l=1}^s \lambda_l \left[ \exp\{\sum_{k=1}^d h_{lk}(y)\theta_k\} / W_l(\theta, G) \right]}$$

Choose the functions  $h_{ik}$  such that the classes

$$(6.1) \quad \begin{aligned} & \{h_{ik}: i = 1, \dots, s - 1, k = 1, \dots, d\} \quad \text{and} \\ & \{w_i(\cdot, \theta)r(\cdot, \theta): \|\theta - \theta_0\|_1 < \delta, i = 1, \dots, s\} \end{aligned}$$

are  $F_{l0}$ -Donsker for some  $\delta > 0$ , for each  $l = 1, \dots, s$ . A sufficient condition for this to hold for the first class of (6.1) is that each  $h_{ik}$  is monotone; see, for example, Van der Vaart and Wellner [(1996), Theorem 2.7.5], which is a natural condition for the vaccine trial application. A sufficient condition for the second class of (6.1) to be Donsker is that each  $w_i r$ , or equivalently each



$r_i$ , is Lipschitz  $\alpha > 1/2$  in  $y$ , uniformly in  $\|\theta - \theta_0\|_1 < \delta$ , and  $\overline{F}_0(|Y|^{2+\delta}) < \infty$  for some  $\delta > 0$ . Here we are using Van der Vaart's (1994) theorem to verify the hypothesis of Ossiander's (1987) bracketing entropy central limit theorem for classes of smooth functions.

**THEOREM 6.1** (Asymptotic properties of the univariate GLR model). *Suppose  $n_i \rightarrow \infty$  and  $\lambda_{ni} \rightarrow \lambda_i > 0$ ,  $i = 1, \dots, s$ , there exists an  $i \in \{1, \dots, s-1\}$  such that the set of functions  $\{h_{i1}, \dots, h_{id}\}$  in (6.1) is linearly independent and the limiting log profile partial likelihood is strictly concave on  $\Theta$ . Further suppose there exists a  $\delta > 0$  such that, for each  $l = 1, \dots, s$ ,*

(A)  $G(h_{ik}^2 \tilde{w}_i(\theta_0)) < \infty$  for each  $i = 1, \dots, s-1$ , and the classes of (6.1) are  $F_{10}$ -Donsker for each  $l = 1, \dots, s$ .

*Then the MLE  $(\hat{\theta}_n, \mathbb{G}_n)$  of  $(\theta_0, G)$  in the GLR model exists uniquely, can be computed by the maximum partial likelihood estimation procedure and is consistent uniformly over  $\mathbf{H} \equiv \{1_{[0,t]}: t \in [0, \infty)\}$ , asymptotically Gaussian in  $R^d \times l^\infty(\mathbf{H})$  and efficient at  $P_0 = P_{(\theta_0, G)}$ .*

The hypothesis  $\lambda_{ni} \rightarrow \lambda_i > 0$ ,  $i = 1, \dots, s$ , states that the fraction of trial participants infected during follow-up does not vanish in the limit for any treatment group. For a two-arm vaccine trial, it simply says that the vaccine is not 100% protective. As stated by Theorem 5 in Gilbert, Lele and Vardi (1999), the hypothesis of a strictly concave limiting log profile partial likelihood holds for all two-sample GLR models. Moreover, an easy extension of the proof of this theorem shows strict concavity in the  $s$ -sample special case with  $h_{1k} = \dots = h_{s-1k}$ ,  $k = 1, \dots, d$ . Thus we have two corollaries.

**COROLLARY 6.1** (Two-sample GLR model). *Suppose  $s = 2$ ,  $n_i \rightarrow \infty$  with  $\lambda_{ni} \rightarrow \lambda_i > 0$ ,  $i = 1, 2$  and (A) holds. Then the conclusion of Theorem 6.1 holds.*

**COROLLARY 6.2** (Special case  $s$ -sample GLR model). *Suppose  $s \geq 2$ ,  $n_i \rightarrow \infty$  with  $\lambda_{ni} \rightarrow \lambda_i > 0$ ,  $i = 1, \dots, s$  and  $h_{1k} = \dots = h_{s-1k}$ ,  $k = 1, \dots, d$ . Then under hypothesis (A), the conclusion of Theorem 6.1 holds.*

Now consider the multivariate GLR model, in which  $\mathbf{Y} = R^k$  and  $\mathbf{H}$  is the collection of indicators of lower-left orthants. Since there are many possibilities for forms of weight functions, we do not state a theorem. Rather, we note that, for a model of interest, the conditions of Theorems 1.1, 2.1, 5.2 and 5.4 can be checked systematically. Essentially, the multivariate situation offers no added difficulties; identifiability and unique estimability in the limit follow under the same conditions, and although the classes (6.1) will contain more functions, they often have entropy bounds under the same hypotheses as for the univariate case. In conclusion, the large sample properties of the MLE established here hold for a broad collection of GLR models.

For two-sample univariate GLR, we explicitly write the limiting variance-covariance of  $(\hat{\theta}_n, \mathbb{G}_n)$ . From (4.6), the limiting variance  $I^{-1}(\theta_0)$  of  $\sqrt{n}(\hat{\theta}_n - \theta_0)$  is

$$(6.2) \quad \left\{ \frac{\lambda_1 \lambda_2}{\int \exp\{h(\theta_0)\} dG} \left[ G\left(\dot{h}^2(\theta_0) * \exp\{h(\theta_0)\} * r(\theta_0)\right) - \frac{G\left[\left(\dot{h}(\theta_0) * \exp\{h(\theta_0)\} * r(\theta_0)\right)^2\right]}{G(\exp\{h(\theta_0)\} * r(\theta_0))} \right] \right\}^{-1}$$

with  $\dot{h}(\theta_0) = \frac{\partial}{\partial \theta} h(\cdot, \theta)|_{\theta_0}$  and  $r^{-1}(y, \theta_0) = \lambda_1(\exp\{h(y, \theta_0)\} / \int \exp\{h(\theta_0)\} dG) + \lambda_2$ . The consistent plug-in estimate of variance is obtained by replacing  $\lambda_i$  and  $(\theta_0, G)$  in (6.2) by  $\lambda_{ni}$  and the MLE  $(\hat{\theta}_n, \mathbb{G}_n)$ . The limiting covariance process  $\{\text{AsymCov}(\mathbb{G}_n(s), \mathbb{G}_n(t)): s, t \in [0, \infty)\}$  and its estimate are obtained by selecting  $h_1 = 1_{[0, s]}$ ,  $h_2 = 1_{[0, t]}$ , and substituting  $w_1, \dot{w}_1$  and  $r$  into (4.9).

EXAMPLE 6.2 (Example 1.3 continued; choice-based sampling in econometrics). In the pure choice-based sampling design of Example 1.3, the  $M$  weight functions in the biased sampling model are  $w_j(y, \theta) = p_\theta(j|y)$ ,  $j = 1, \dots, M$ ,  $y \in \mathbf{Y}$ . In this application, the sample space  $\mathbf{Y} \subset R^d$ , and  $\mathbf{H}$  is the set of indicators of lower-left orthants times the indicators of  $\mathbf{Y}$ . Since every weight function depends on  $\theta$ , Theorem 1.1 cannot be used to establish identifiability of the model. However, Theorem 3 of Gilbert, Lele and Vardi (1999) gives a sufficient condition for identifiability in this case. It yields, for instance, that the semiparametric biased sampling model is identifiable when the parametric model  $p_\theta$  is a baseline category logit model [Cox and Snell (1989); Fienberg (1980)] given by

$$p_\theta(j|y) = \frac{\exp(\beta_j^T y)}{\sum_{j=1}^M \exp(\beta_j^T y)}, \quad j = 1, \dots, M,$$

or an adjacent categories linear logit model [Agresti (1984); Fienberg (1980)] given by

$$p_\theta(j|y) = \frac{\exp(\alpha_j + j\beta^T y)}{\sum_{j=1}^M \exp(\alpha_{j'} + j'\beta^T y)}, \quad j = 1, \dots, M.$$

When  $p_\theta$  is the full multinomial logistic regression model (1.7), the biased sampling model can also be shown to be identifiable, arguing as in the proof of Theorem 3 in Gilbert, Lele and Vardi (1999), utilizing the fact that the aggregate shares are known. Moreover, the maximum partial likelihood estimation procedure applies because the aggregate shares are known.

For this problem,  $\dot{w}_{jk}(y, \theta)/w_j(y, \theta) = \dot{l}_{\theta_k}(j|y)$ , the  $k$ th component of the score for the parametric model  $p_\theta(j|y)$ . Suppose there exists a  $\delta > 0$  such that  $\sup_{\|\theta - \theta_0\|_1 < \delta} G(|\dot{l}_{\theta_k}(j|Y)|^2 p_{\theta_0}(j|Y)) < \infty$ , which usually holds, for instance, if the covariates are assumed bounded. This immediately implies condition (vii)  $\sup_{\|\theta - \theta_0\|_1 < \delta} G(|\dot{w}_{kl}(\theta)/w_k(\theta)|^2 \tilde{w}_k(\theta_0)) < \infty$  of Theorem 5.3, and consistency

of the MLE  $(\hat{\theta}_n, \mathbb{G}_n)$  also follows if each  $\dot{l}_{\theta_k}(j|\cdot)$  is continuous in  $\theta$  at  $\theta_0$  and  $l_{1\text{pro}}(\theta)$  is strictly concave on  $\Theta$ . As described in Remark 5.1,  $l_{1\text{pro}}(\theta)$  is strictly concave on  $\Theta$  and (iii) holds if the information  $I(\theta)$  for  $\theta$  in the biased sampling model is positive definite for all  $\theta \in \Theta$ . For most parametric models  $p_\theta$  defining the biasing functions, this is the case, which can be directly verified by checking condition (I).

The uniform Lipschitz condition (viii) of Theorem 5.3 takes the form

$$|\dot{l}_{\theta_{1k}}(j|y) - \dot{l}_{\theta_{2k}}(j|y)| \leq K\|\theta_1 - \theta_2\|_1^\alpha$$

for all  $\theta_1, \theta_2 \in \{\theta \in \Theta: \|\theta - \theta_0\|_1 < \delta\}$ . For the case in which  $p_\theta$  has form (1.7),  $\theta = [\alpha_1, \dots, \alpha_M, \beta_1, \dots, \beta_M]^T$ , and straightforward calculation shows that, for  $k = 1, \dots, M$  (i.e.,  $\theta_1 = \alpha_1, \dots, \theta_M = \alpha_M$ ),

$$\begin{aligned} \left[ \dot{l}_{\theta_{1k}}(j|y) - \dot{l}_{\theta_{2k}}(j|y) \right] &= [p_{\theta_2}(k|y) - p_{\theta_1}(k|y)] \\ (6.3) \qquad \qquad \qquad &= \left[ \frac{\exp(\alpha_{2k} + \beta_{2k}^T y)}{\sum_{j'=1}^M \exp(\alpha_{2j'} + \beta_{2j'}^T y)} - \frac{\exp(\alpha_{1k} + \beta_{1k}^T y)}{\sum_{j'=1}^M \exp(\alpha_{1j'} + \beta_{1j'}^T y)} \right]. \end{aligned}$$

For  $k = M + 1, \dots, 2M$  (i.e.,  $\theta_{M+1} = \beta_1, \dots, \theta_{2M} = \beta_M$ ),

$$\left[ \dot{l}_{\theta_{1k}}(j|y) - \dot{l}_{\theta_{2k}}(j|y) \right] = y[p_{\theta_2}(k - M|y) - p_{\theta_1}(k - M|y)].$$

Thus, the uniform Lipschitz condition will hold if the covariate space  $\mathbf{Y}$  is bounded and, for each  $k$ , (6.3) is bounded by a constant times  $\|\theta_1 - \theta_2\|_1^\alpha$  for some  $\alpha > d/2$ .

Then asymptotic normality and efficiency of  $(\hat{\theta}_n, \mathbb{G}_n)$  follow as long as the two classes

$$\begin{aligned} &\{\dot{l}_{\theta_0 k}(j|\cdot): j = 1, \dots, M, k = 1, \dots, d\}, \\ &\left\{ \frac{p_\theta(j|\cdot)}{\sum_{j'=1}^M \lambda_{j'} [p_\theta(j'|\cdot)/G(p_\theta(j'|\cdot))]} : \|\theta - \theta_0\|_1 < \delta, j = 1, \dots, M \right\}, \end{aligned}$$

are  $F_{l0}$ -Donsker for each  $l = 1, \dots, s$ , for some  $\delta > 0$ . The limiting covariance structure of  $(\hat{\theta}_n, \mathbb{G}_n)$  can be explicitly written down by substituting the chosen model  $p_{\theta_0}$  into the information bound formulae.

**7. Proofs.**

PROOF OF THEOREM 4.1. The efficient score function for  $\theta$  in  $\mathbf{P}$  can be calculated by identifying a tangent  $a' \in L_2^0(G)$  satisfying  $E[(\dot{l}_\theta - \dot{l}_g a')(\dot{l}_g a')] = 0$  for all  $a \in L_2^0(G)$ , and verifying that the tangent space of  $G, \dot{\mathbf{P}}_2$ , equals  $\mathbf{R}(\dot{l}_g)$ ,

the closure of the range of the score for  $g$ . Using (4.1) and (4.2), we calculate

$$\begin{aligned}
 0 &= E\left[\{\dot{l}_\theta(I, Y) - \dot{l}_g a'(I, Y)\}\{\dot{l}_g a(I, Y)\}\right] \\
 &= E\left[\left[\frac{\dot{w}_I(Y, \theta)}{w_I(Y, \theta)} - a'(Y)\right]a(Y)\right] \\
 (7.1) \quad &- E\left[E\left[\frac{\dot{w}_I(Y, \theta)}{w_I(Y, \theta)} - a'(Y)|I\right]E[a(Y)|I]\right] \\
 &= \int \left\{ \sum_{i=1}^s \lambda_i \left[ \frac{\dot{w}_i(y, \theta)}{w_i(y, \theta)} - a'(y) \right. \right. \\
 &\quad \left. \left. - E\left[\frac{\dot{w}_i(Y, \theta)}{w_i(Y, \theta)} - a'(Y)|I = i\right] \right] \tilde{w}_i(y, \theta) \right\} a(y) dG(y).
 \end{aligned}$$

Since this equality holds for all  $a \in L_2^0(G)$ , it follows that the quantity in braces must vanish almost everywhere in  $L_2^0(G)$ , which leads to a Fredholm integral equation of the second kind:

$$(7.2) \quad a'(y) - \int a'(z)K(y, z) dG(z) = \phi^*(y),$$

with kernel  $K$  and function  $\phi^*$  given by  $K(y, z) = \sum_{i=1}^s r_i(y, \theta)\tilde{w}_i(z, \theta)$  and  $\phi^*(y) = \sum_{i=1}^s r_i(y, \theta)\dot{l}_\theta(i, y)$ .

We suppress the dependence of  $a', K$  and  $\phi^*$  on  $\theta$  and  $G$ . As in Tricomi (1957), since the kernel is Pincherle–Goursat (i.e., it factors in  $y$  and  $z$ ), the Fredholm equation (7.2) can be equivalently expressed as an algebraic system of  $s$  equations in  $s$  unknowns  $A\underline{z} = B$ , where  $\underline{z} = (z_1, \dots, z_s)^T$  are unknowns,  $A = I_s - (a_{ij}) = I_s - G(r\underline{\tilde{w}}^T)$  and  $B = [b_1^T, \dots, b_s^T]^T$ , with

$$\begin{aligned}
 (7.3) \quad a_{ij} &= \int r_j(y)\tilde{w}_i(y) dG(y) = \lambda_j G(r\tilde{w}_i\tilde{w}_j), \\
 b_j &= \int \phi^*(y)\tilde{w}_j(y) dG(y) = \sum_{i=1}^s \lambda_i G(\dot{l}_\theta(i, Y)r\tilde{w}_i\tilde{w}_j).
 \end{aligned}$$

It is easily seen that  $A = M\underline{\lambda}$ , where  $M$  is the matrix defined in (3.1), so that  $A$  has rank  $s - 1$ . Letting  $\overline{M}^-$  be a generalized inverse of  $M$ , we solve  $\underline{z} = \underline{\lambda}^{-1}\overline{M}^-B$ , which implies

$$a'(y) = \phi^*(y) + B^T(M^-)^T\underline{\lambda}^{-1}r = \phi^*(y) + G(\phi^*\underline{\tilde{w}}^T)(M^-)^T r\underline{\tilde{w}}.$$

Set  $\dot{l}_\theta^* = \dot{l}_\theta - \dot{l}_g a'$ , which is unique since  $(\dot{l}_g a')(x) = a'(y) - G(a'\tilde{w}_i)$  is independent of the choice of  $M^-$ . This follows by Lemma 5.2(iv) of GVW, since  $r\underline{\tilde{w}} - G(r\underline{\tilde{w}}\tilde{w}_i) \in \text{Range}(M)$ . Since  $G(a')$  is not necessarily equal to zero, we center by taking  $a_* = a' - G(a') \in \underline{L}_2^0(G)$ , and set  $\dot{l}_\theta^* = \dot{l}_\theta - \dot{l}_g a_* = \dot{l}_\theta - \dot{l}_g a'$ , which is unique and orthogonal to  $\mathbf{R}(\dot{l}_g)$ . A slight adaptation of Theorem 3 of

BKRW (page 123) proves that if the graph  $\mathbf{G}^*(\theta)$  is connected for all  $\theta$  in some neighborhood of  $\theta_0$ , then  $\dot{\mathbf{P}}_2 = \{a - E(a|I) : a \in L_2^0(G)\} = \overline{\mathbf{R}(\dot{l}_g)}$ , so that  $\dot{l}_\theta^*$  is orthogonal to  $\dot{\mathbf{P}}_2$  and is therefore the efficient score function.

Nonsingularity of  $I(\theta_0)$  is shown in the proof of Theorem 5.3. Straightforward calculation, aided with the identity  $\sum_{i=1}^s \lambda_i G(\tilde{w}_i r) = 1$ , yields the form (4.6) of the inverse information when  $s = 2$ .

PROOF OF THEOREM 4.2. Consider the parameter  $G$  as a functional of the model:  $\nu(P_{(\theta,G)}) = \chi(G) = G$ . Let  $\dot{\chi}$  denote the pathwise derivative of  $\chi$ , and  $\dot{\chi}^T$  its adjoint. Let  $(\dot{l}_g^T \dot{l}_g)^*$  be the information bound operator for estimation of  $G$ , viewed as a map from  $L_2(G)$  and  $L_2(G)$ . We apply Corollary 3 of Begun, Hall, Huang and Wellner (1983), proved in BKRW (pages 218–220), which states that if  $\nu$  is pathwise differentiable at  $P_0$ ,  $\dot{\mathbf{P}} = \mathbf{R}(\dot{l}_g) + \mathbf{R}(\dot{l}_\theta)$ ,  $\mathbf{R}(\dot{\chi}^T) \subset \mathbf{R}(\dot{l}_g^T \dot{l}_g)$ , and  $(\dot{l}_g^T \dot{l}_g)^{\star-1}$  and  $I^{-1}(\theta)$  exist, then the inverse information covariance functional for  $\nu$  equals

$$(7.4) \quad \begin{aligned} I_\nu^{-1}(h_1, h_2) &= E \left[ \dot{\chi}^T \pi_{h_1} (\dot{l}_g^T \dot{l}_g)^{\star-1} \dot{\chi}^T \pi_{h_2} \right] \\ &\quad + E \left[ \dot{\chi}^T \pi_{h_1} \alpha_\star \right]^T I^{-1}(\theta) E \left[ \dot{\chi}^T \pi_{h_2} \alpha_\star \right], \end{aligned}$$

where  $\alpha_\star$  is the unique tangent in  $L_2^0(G)$  satisfying  $E[(\dot{l}_\theta - \dot{l}_g \alpha_\star)(\dot{l}_g a)] = 0$  for all  $a \in L_2^0(G)$ .

Since  $\chi$  is just the identity, it is straightforward to show [see BKRW, Example 1, page 191 and BKRW, panel (5), page 58] that  $\dot{\chi}^T \pi_h(y) = h(y) - G(h)$  for a projection map  $\pi_h$ ,  $h \in \mathbf{H}$ . The operator  $(\dot{l}_g^T \dot{l}_g)^{\star-1}$  is calculated by GVW (pages 1089–1091), and requires that  $r(y, \theta_0)$  and  $r^{-1}(y, \theta_0)$  are bounded on  $\mathbf{Y}$ . Substituting it and  $\dot{\chi}^T \pi_h(y) = h(y) - G(h)$  into (7.4) and calculation shows that  $I_\nu^{-1}(h_1, h_2)$  equals expression (4.8) with  $\alpha_\star$  as in the proof of Theorem 4.1.  $\square$

PROOF OF THEOREM 5.1. For fixed  $\theta$ , let  $\underline{v}_n(\theta) = (V_{n1}(\theta), \dots, V_{ns-1}(\theta), 1)^T$  be the solution of the system (1.12). Since  $\theta$  is fixed,  $\underline{v}_n(\theta) \rightarrow_{\text{a.s.}} \underline{V}(\theta)$  by Proposition 2.1 of GVW [page 1081]. Since each  $\lambda_{ni}^{-1} r_{ni}(\cdot, \theta)$  is a bounded function of  $y$ , if the processes

$$(7.5) \quad \left\{ C_{ni}(v) \equiv \frac{1}{n_i} \sum_{j=1}^{n_i} \log \left\{ \frac{w_{ij}(\theta) v_i^{-1}}{\sum_{k=1}^s \lambda_{nk} w_{kj}(\theta) v_k^{-1}} \right\} : v = (v_1, \dots, v_s)^T \in V \right\},$$

$$V \equiv [0, \infty)^{s-1} \times 1$$

converge uniformly on  $V$ , then the log profile partial likelihood  $l_{n1\text{pro}}$  satisfies, for each fixed  $\theta \in \Theta$

$$\frac{1}{n} l_{n1\text{pro}}(\theta|\underline{x}) \rightarrow_{\text{a.s.}} \sum_{i=1}^s \lambda_i \int \log \left\{ \frac{w_i(y, \theta) V_i^{-1}(\theta)}{\sum_{k=1}^s \lambda_k w_k(y, \theta) V_k^{-1}(\theta)} \right\} \tilde{w}_i(y, \theta_0) dG(y)$$

$$= l_{1\text{ pro}}(\theta).$$

Since  $l_{1\text{ pro}}(\theta)$  and  $l_{n1\text{ pro}}(\theta|\underline{x})$  for  $n > N$  are strictly concave on  $\Theta$ , this implies  $\hat{\theta}_n \rightarrow_{\text{a.s.}} \theta_0$  as  $n \rightarrow \infty$ , applying Theorem II.1 in Appendix II of Andersen and Gill (1982).

Uniform convergence of the processes in (7.5) can be established by the bracketing entropy Glivenko–Cantelli Theorem 2.4.1 of Van der Vaart and Wellner (1996).  $\square$

To prove consistency of  $\mathbb{G}_n$ , some additional notation is needed. Where possible, it matches the notation used by GVW. Define sets of  $s - 1$  equations, which correspond to (1.12), by

$$(7.6) \quad 1 = \mathbb{H}_{ni}(\mathbb{V}_1, \dots, \mathbb{V}_{s-1}, 1, \theta), \quad 1 = H_{ni}(V_1, \dots, V_{s-1}, 1, \theta)$$

for  $i = 1, \dots, s - 1$ ,  $\theta \in \Theta$ , where

$$\mathbb{H}_{ni}(u_1, \dots, u_s, \theta) \equiv \frac{1}{u_i} \int \frac{w_i(y, \theta)}{\sum_{j=1}^s (\lambda_{nj} w_j(y, \theta)) / u_j} d\mathbb{F}_n(y)$$

and  $H_{ni}$  equals  $\mathbb{H}_{ni}$  with  $\mathbb{F}_n$  replaced by  $\overline{F}_{n0}$ . Reparametrize by setting  $z_j = \log\{\lambda_{nj}/u_j\}$ , and define  $\mathbb{K}_{ni}(z_1, \dots, z_s, \theta) = \mathbb{H}_{ni}(\lambda_{n1}e^{-z_1}, \dots, \lambda_{ns}e^{-z_s}, \theta) - \lambda_{ni}$  and  $K_{ni}(z_1, \dots, z_s, \theta) = H_{ni}(\lambda_{n1}e^{-z_1}, \dots, \lambda_{ns}e^{-z_s}, \theta) - \lambda_{ni}$ . Then (7.6) become

$$0 = \mathbb{K}_{ni}(\mathbb{Z}_1, \dots, \mathbb{Z}_{s-1}, \log\{\lambda_{ns}\}\theta), \quad 0 = K_{ni}(Z_1, \dots, Z_{s-1}, \log\{\lambda_{ns}\}, \theta)$$

for  $i = 1, \dots, s - 1$ ,  $\theta \in \Theta$ . Set  $\underline{\mathbb{H}}_n = (\mathbb{H}_{n1}, \dots, \mathbb{H}_{ns})^T$ , and similarly defined  $\underline{H}_n$ ,  $\underline{\mathbb{K}}_n$  and  $\underline{K}_n$ . Observe that  $\underline{\mathbb{K}}_n(z, \theta) = \nabla \mathbb{D}_n(z, \theta)$ , where  $\mathbb{D}_n: R^s \times \Theta \rightarrow R^1$  is given by

$$\mathbb{D}_n(z, \theta) = \int \log \left[ \sum_{i=1}^s e^{z_i} w_i(y, \theta) \right] d\mathbb{F}_n(y) - \sum_{i=1}^s \lambda_{ni} z_i.$$

Similarly, we can write  $\underline{K}_n(z, \theta) = \nabla D_n(z, \theta)$ , where  $D_n$  equals  $\mathbb{D}_n$  with  $\mathbb{F}_n$  replaced by  $\overline{F}_{n0}$ .

As proved in GVW (pages 1079–1080), if the graph  $\mathbf{G}^*(\theta)$  is connected, then  $\mathbb{D}_n(z, \theta)$  and  $D_n(z, \theta)$  are convex functions of  $\underline{z}$  ( $\theta$  fixed). Thus the reparametrization allows us to use convexity theory to conclude consistency.

For an  $s$ -vector  $\underline{u} \in R^s$ , define

$$\begin{aligned} \tilde{w}(\underline{u}) &\equiv \text{diag}(u^{-1})\underline{w}, & r(\underline{u}) &\equiv (\underline{\lambda}^T \tilde{w}(\underline{u}))^{-1}, \\ r_n(\underline{u}) &\equiv (\underline{\lambda}_n^T \tilde{w}(\underline{u}))^{-1}, & r_{ni}(\underline{u}) &\equiv (\lambda_{ni} \tilde{w}_i r_n(\underline{u})). \end{aligned}$$

Therefore,  $\tilde{w}$ ,  $r$ ,  $r_n$  and  $r_{ni}$  are given by

$$\begin{aligned} \tilde{w} &= \tilde{w}(\underline{W}) = \text{diag}(\underline{W}^{-1})\underline{w}, & r &\equiv (\underline{\lambda}^T \tilde{w})^{-1}, \\ r_n &\equiv (\underline{\lambda}_n^T \tilde{w})^{-1}, & r_{ni} &\equiv (\lambda_{ni} \tilde{w}_i r_n). \end{aligned}$$

We can rewrite the functions  $\mathbb{H}_n$ ,  $\underline{H}_n$ ,  $\mathbb{D}_n$  and  $D_n$  in this notation:

$$\begin{aligned} \mathbb{H}_n(\underline{u}, \theta) &= \int r_n(\underline{u}, \theta) \underline{w}(\underline{u}, \theta) d\mathbb{F}_n, \\ \underline{H}_n(\underline{u}, \theta) &= \int r_n(\underline{u}, \theta) \underline{w}(\underline{u}, \theta) d\overline{F}_{n0}, \\ \mathbb{D}_n(\underline{z}, \theta) &= \int \log \left[ \sum_{i=1}^s e^{z_i} w_i(\theta) \right] d\mathbb{F}_n - \sum_{i=1}^s \lambda_{ni} z_i, \\ D_n(\underline{z}, \theta) &= \int \log \left[ \sum_{i=1}^s e^{z_i} w_i(\theta) \right] d\overline{F}_{n0} - \sum_{i=1}^s \lambda_{ni} z_i. \end{aligned}$$

Introduce functions  $\underline{H}$  and  $\underline{D}$  equal to  $\underline{H}_n$  and  $\underline{D}_n$ , respectively, with  $\lambda_{ni}$ ,  $r_n$  and  $\overline{F}_{n0}$  replaced with  $\lambda_i$ ,  $r$  and  $\overline{F}_0$ . Finally, define  $\tilde{\mathbb{D}}_n(\underline{z}, \theta) = \mathbb{D}_n(\underline{z}, \theta) - \mathbb{D}_n(\underline{Z}, \theta)$  and  $\tilde{D}(\underline{z}) = D(\underline{z}) - D(\underline{Z})$ , where  $\underline{Z} \equiv (Z_1, \dots, Z_s)^T$ ,  $Z_i \equiv \log\{\lambda_i/V_{i0}\}$ . Since  $\mathbb{D}_n(\underline{z}, \theta)$  and  $D(\underline{z}, \theta)$  are convex functions of  $\underline{z}$  for each  $\theta \in \Theta$ , so are  $\tilde{\mathbb{D}}_n$  and  $\tilde{D}$ .

Using this notation, we establish consistency, which depends on the uniqueness of  $\underline{V}_n(\hat{\theta}_n)$  and  $\underline{V}(\theta_0)$  as solutions of the systems of equations (7.6), together with convexity of  $\mathbb{D}_n(\cdot, \hat{\theta}_n)$  and  $\tilde{D}(\cdot, \theta_0)$ .

The following lemma generalizes Lemma 5.3 in GVW, allowing the function  $\tilde{\mathbb{D}}_n$  to depend on  $\theta$ .

LEMMA 7.1. *Assume hypotheses (ii) of Theorem 5.1 and (iv) of Proposition 5.1. Then*

$$(7.7) \quad \tilde{\mathbb{D}}_n(\underline{z}, \hat{\theta}_n) \rightarrow_{\text{a.s.}} \tilde{D}(\underline{z}, \theta_0)$$

as  $n \rightarrow \infty$  with  $\lambda_{ni} \rightarrow \lambda_i > 0$  for each fixed  $\underline{z} \in R^s$ . Since  $\tilde{\mathbb{D}}_n$  is convex, it also holds that

$$(7.8) \quad \sup_{\underline{z} \in C} |\tilde{\mathbb{D}}_n(\underline{z}, \hat{\theta}_n) - \tilde{D}(\underline{z}, \theta_0)| \rightarrow_{\text{a.s.}} 0$$

as  $n \rightarrow \infty$  with  $\lambda_{ni} \rightarrow \lambda_i > 0$  for any compact subset  $C \subset R^s$ .

PROOF OF LEMMA 7.1. Define

$$q(\underline{z}, y, \theta) \equiv \log \left\{ \frac{\sum_{i=1}^s e^{z_i} w_i(y, \theta)}{\sum_{i=1}^s e^{Z_i} w_i(y, \theta)} \right\}.$$

Observe that  $q(\underline{z}, y, \theta)$  is a bounded function of  $y$  and  $\theta$ . Let  $\Theta_0$  be a neighborhood of  $\theta_0$  for which (iv) of Proposition 5.1 holds. Towards showing (7.7), we have

$$\tilde{\mathbb{D}}_n(\underline{z}, \hat{\theta}_n) = \sum_{j=1}^s \lambda_{nj} \mathbb{F}_{nj} \left[ q(\underline{z}, Y, \hat{\theta}_n) \right] - \sum_{i=1}^s \lambda_{ni} (z_i - Z_i).$$

Expand  $q$  about  $\theta_0$ , giving, for fixed  $\underline{z}$ ,

$$\begin{aligned} \mathbb{F}_{nj}[q(\underline{z}, Y, \hat{\theta}_n)] &= \mathbb{F}_{nj}[q(\underline{z}, Y, \theta_0)] \\ &\quad + (\hat{\theta}_n - \theta_0)^T * \mathbb{F}_{nj} \left[ \frac{\sum_{i=1}^s e^{z_i} \dot{w}_i(Y, \theta^*)}{\sum_{i=1}^s e^{z_i} w_i(Y, \theta^*)} - \frac{\sum_{i=1}^s e^{Z_i} \dot{w}_i(Y, \theta^*)}{\sum_{i=1}^s e^{Z_i} w_i(Y, \theta^*)} \right] \end{aligned}$$

for a  $\theta^*$  between  $\theta_0$  and  $\hat{\theta}_n$ . The first piece converges almost surely to  $\int q(\underline{z}, y, \theta_0) \bar{w}_j(y, \theta_0) dG(y)$  by the strong law of large numbers (as in Lemma 5.3 of GVW), while the second piece is  $o_p(1)$  because  $\hat{\theta}_n - \theta_0 \rightarrow_p \underline{0}$  and the  $\mathbb{F}_{nj}[\cdot]$  term is  $O_p(1)$  for each  $j = 1, \dots, s$  and  $k = 1, \dots, d$ . Here continuity of each  $w_i$  and  $\dot{w}_i/w_i$  in  $\theta$  at  $\theta_0$  and hypothesis (iv) of Proposition 5.1 are used. Thus (7.7) holds.

Since  $\tilde{\mathbb{D}}_n$  is convex, (7.8) follows from (7.7), using Theorem II.1 in Appendix II of Andersen and Gill (1982) [also see Exercise 3.2.4 of Van der Vaart and Wellner (1996)].  $\square$

**PROOF OF PROPOSITION 5.1.** Our proof closely follows the proof of Proposition 2.1 in GVW. Choose  $N$  so that, for  $n > N$ ,  $\hat{\theta}_n$  is in a neighborhood  $\Theta_0$  of  $\theta_0$  for which (iv) of Proposition 5.1 holds. Fix  $n > N$ . Since the graph  $\mathbf{G}^*(\hat{\theta}_n)$  is connected, Lemma 1 of Vardi [(1985), page 197] implies that  $\underline{V}(\hat{\theta}_n)$  is the unique solution of the right equation system in (7.6). Therefore,  $\underline{Z}_n = (Z_{n1}, \dots, Z_{ns})^T$ , with  $Z_{ni} = \log\{\lambda_i/V_i(\hat{\theta}_n)\}$ , is the unique minimizer of  $D(\underline{z}, \hat{\theta}_n)$ , and thus also of  $\tilde{D}(\underline{z}, \hat{\theta}_n)$ , subject to  $z_s = \log\{\lambda_s\}$ . Let  $\underline{Z}_n$  be the corresponding minimizer of  $D_n(\underline{z}, \hat{\theta}_n)$ , which also minimizes  $\tilde{D}_n(\underline{z}, \hat{\theta}_n)$  subject to  $z_s = \log\{\lambda_{ns}\}$ . Also let  $\underline{Z}$  be the unique minimizer of  $\tilde{D}(\underline{z}, \theta_0)$ . Then for any compact set  $C \in R^s$  with  $\underline{Z}$  in the interior of  $C$ , it follows from (7.7) and (7.8) of Lemma 7.1, and the definition of  $\tilde{D}(\underline{z}, \theta)$  and  $\underline{Z}$ , that  $\inf_{\underline{z} \in \partial C} \tilde{\mathbb{D}}_n(\underline{z}, \hat{\theta}_n) \rightarrow_{a.s.} \inf_{\underline{z} \in \partial C} \tilde{D}_n(\underline{z}, \theta_0) > 0$ , while  $\tilde{\mathbb{D}}_n(\underline{Z}_n, \hat{\theta}_n) \rightarrow_{a.s.} \tilde{D}(\underline{Z}, \theta_0) = 0$ . The convexity of  $\tilde{\mathbb{D}}_n(\underline{z}, \theta)$  in  $\underline{z}$  for all  $\theta \in \Theta_0$  implies that  $\underline{Z}_n \in C$  for all  $n$  greater than some  $N$  with probability 1. Since  $C$  can be made arbitrarily small, a simple  $\varepsilon - \delta$  argument yields  $\underline{Z}_n \rightarrow_{a.s.} \underline{Z}$ . Since  $\underline{Z}_n = (Z_{n1}, \dots, Z_{ns})^T$  with  $Z_{ni} = \log\{\lambda_{ni}/V_{ni}(\hat{\theta}_n)\}$ , where  $\underline{V}_n(\hat{\theta}_n) = (V_{n1}(\hat{\theta}_n), \dots, V_{ns-1}(\hat{\theta}_n), 1)^T$  is the solution of the left equation system in (7.6), it follows that  $\underline{V}_n(\hat{\theta}_n) \rightarrow_{a.s.} \underline{V}(\theta_0)$ . The proof that  $\underline{W}_n \rightarrow_{a.s.} \underline{W}(\theta_0)$  is similar and is omitted.  $\square$

**PROOF OF THEOREM 5.2.** This proof closely follows the proof of Theorem 2.1 in GVW. Fix  $h \in \mathbf{H}$ , with  $\mathbf{H}$  defined in (5.3), which satisfies  $G(h) < \infty$ . Write

$$\begin{aligned} (7.9) \quad & \sup_{h \in \mathbf{H}} |\mathbb{F}_n(hr_n(\underline{V}_n(\hat{\theta}_n))) - \bar{F}_0(hr(\theta_0))| \\ & \leq \sup_{h \in \mathbf{H}} \left| \mathbb{F}_n \left( h \left[ r_n(\underline{V}_n(\hat{\theta}_n)) - r(\hat{\theta}_n) \right] \right) \right| \\ & \quad + \sup_{h \in \mathbf{H}} |\mathbb{F}_n(hr(\hat{\theta}_n)) - \bar{F}_0(hr(\theta_0))| \end{aligned}$$



$$(7.10) \quad \leq \left\| \frac{r_n(\underline{V}_n(\hat{\theta}_n))}{r(\hat{\theta}_n)} - 1 \right\|_\infty |\mathbb{F}_n(h_e r(\hat{\theta}_n))| + \sup_{h \in \mathbf{H}} |\mathbb{F}_n(hr(\hat{\theta}_n)) - \bar{F}_0(hr(\theta_0))|.$$

Since  $\hat{\theta}_n \rightarrow_{\text{a.s.}} \theta_0$  and  $\underline{V}_n(\hat{\theta}_n) \rightarrow_{\text{a.s.}} \underline{V}(\theta_0)$  by Theorem 5.1 and Proposition 5.1, and  $r_n(\underline{u})/r$  is continuous and bounded in  $y$  and  $\theta$ , the  $\|\cdot\|_\infty$  term converges almost surely to zero. The term  $|\mathbb{F}_n(h_e r(\hat{\theta}_n))|$  is stochastically bounded. Next, we show that

$$(7.11) \quad \sup_{h \in \mathbf{H}} |\mathbb{F}_n(hr_n(\hat{\theta}_n)) - \bar{F}_0(hr(\theta_0))| \rightarrow_{\text{a.s.}} 0.$$

Expanding about  $\theta_0$  by differentiating  $r_n(y, \theta)$  with respect to  $\theta$ , (7.11) is bounded above by

$$(7.12) \sup_{h \in \mathbf{H}} |\mathbb{F}_n(hr_n(\theta_0)) - \bar{F}_0(hr(\theta_0))| + \sum_{i=1}^s \lambda_{ni} \sup_{h \in \mathbf{H}} |(\hat{\theta}_n - \theta_0)^T * \int h(y) \left( - \sum_{k=1}^s \lambda_{nk} \left[ \frac{\dot{w}_k(y, \theta^*)}{w_k(y, \theta^*)} - \int \frac{\dot{w}_k(z, \theta^*)}{w_k(z, \theta^*)} \tilde{w}_k(z, \theta^*, \mathbb{F}_n) d\mathbb{F}_n(z) \right] \times \tilde{w}_k(y, \theta^*, \mathbb{F}_n) r_n^2(y, \theta^*) \right) d\mathbb{F}_n(y) |,$$

where  $\tilde{w}_k(z, \theta, \mathbb{F}_n) \equiv w_k(z, \theta)/W_k(\theta, \mathbb{F}_n)$ . Exactly as in GVW, piece (7.12) converges to zero almost surely by Pollard’s Glivenko–Cantelli theorem [see Dudley (1984), Theorems 11.1.2 and 11.1.6] for  $\mathbf{H}$  as in (5.3). By (iv), (vi) and (3.2), and since  $h \in L_1(G)$ , each supremum in (7.13) is  $o_p(1)$ , so that (7.11) holds. It follows from (7.11) and (7.10) that piece (7.9) converges almost surely to zero. Then

$$|\mathbb{G}_n(h) - G(h)| = \left| \frac{\mathbb{F}_n(hr_n(\underline{V}_n(\hat{\theta}_n)))}{\mathbb{F}_n(r_n(\underline{V}_n(\hat{\theta}_n)))} - \frac{\bar{F}_0(hr(\theta_0))}{\bar{F}_0(r(\theta_0))} \right| \leq \frac{|\mathbb{F}_n(hr_n(\underline{V}_n(\hat{\theta}_n))) - \bar{F}_0(hr(\theta_0))|}{\mathbb{F}_n(r_n(\underline{V}_n(\hat{\theta}_n)))} + \frac{\bar{F}_0(hr(\theta_0))}{\bar{F}_0(r(\theta_0))} \frac{|\mathbb{F}_n(r_n(\underline{V}_n(\hat{\theta}_n))) - \bar{F}_0(r(\theta_0))|}{\mathbb{F}_n(r_n(\underline{V}_n(\hat{\theta}_n)))},$$

so (5.4) follows from (7.11) with  $h = 1$ .  $\square$

**PROOF OF THEOREM 5.3.** We apply the Master theorem; see, for example, BKRW (Theorem 1, page 312). First we calculate the matrix  $\Psi_0$ . Differentia-

tion of  $P_0\psi_0(X)$  gives

$$(7.14) \quad \begin{aligned} \dot{\Psi}_0 = P_0 \nabla_{\theta} \left[ \frac{\dot{w}_I(Y, \theta)}{w_I(Y, \theta)} - \phi(Y, \theta) \right] \Big|_{\theta_0} \\ + \nabla_{\theta} \underline{d}(\theta) \Big|_{\theta_0} P_0 [-\underline{e}_I + \underline{r}(Y, \theta_0)] + \underline{d}(\theta_0) P_0 \nabla_{\theta} [\underline{r}(Y, \theta)] \Big|_{\theta_0}. \end{aligned}$$

The term  $P_0[-\underline{e}_I + \underline{r}(Y, \theta_0)] = 0$ , so that it is not necessary to compute  $\nabla_{\theta} \underline{d}(\theta) \Big|_{\theta_0}$ . Direct differentiation using the chain rule and the property  $M^-MM^- = M^-$  of a  $\{1, 2\}$ -generalized inverse shows that the rightmost term of (7.14) is also the zero-vector. Thus  $\dot{\Psi}_0$  equals the first term on the right-hand side of (7.14). Then, differentiation (again using the chain rule) and tedious calculation (which can be readily checked for the case  $s = 2$ ) shows that  $\dot{\Psi}_0$  equals  $-I(\theta_0)$  as displayed in (4.4).

Next we show invertibility of  $\dot{\Psi}_0 = -I(\theta_0)$ . Consider  $I(\theta_0)$  expressed as an expected conditional variance as in Remark 4.3. For  $b \in R^d$ , evidently  $b^T I(\theta_0)b \geq 0$ , with equality if and only if

$$\text{Var} \left( b^T \left[ \frac{\dot{w}_I(Y, \theta_0)}{w_I(Y, \theta_0)} - \alpha_*(Y, \theta_0) \right] \Big| I = i \right) = 0$$

for each  $i = 1, \dots, s$ , which holds if and only if  $b^T [(\dot{w}_i(y, \theta_0)/w_i(y, \theta_0)) - \alpha_*(y, \theta_0)]$  is constant a.s. with respect to the law  $F_{i0}$ , for each  $i = 1, \dots, s$ . Thus (I) implies  $\dot{\Psi}_0$  and  $I(\theta_0)$  are nonsingular.

Direct calculation shows that  $P_0[\psi_{\theta_0}(X)\psi_{\theta_0}(X)^T]$ , where the influence function  $\psi_{\theta_0}$  is as defined in (5.11), equals  $I(\theta_0)$ . Therefore,  $\Sigma = (-I^{-1}(\theta_0))I(\theta_0)(-I^{-1}(\theta_0))^T = I^{-1}(\theta_0)$ .

Next we verify (5.12), that is, that the linearization condition of the Master theorem [(GM2) of BKRW Theorem 1, page 312] holds. Writing  $\mathbb{P}_n[\psi_n(X, \theta, \underline{V}_n(\theta, \mathbb{P}_n)) - \psi(X, \theta, \underline{V}(\theta, P))]$  as  $(\mathbb{P}_n - P)[\cdot] + P[\cdot]$ , the  $(\mathbb{P}_n - P)[\cdot]$  term can straightforwardly be shown to be  $o_p(n^{-1/2})$ , so it suffices to show that

$$(7.15) \quad \begin{aligned} & P \left[ \psi_n(X, \theta, \underline{V}_n(\theta, \mathbb{P}_n)) - \psi(X, \theta, \underline{V}(\theta, P)) \right] \\ & = P \left[ \psi_n(X, \theta, \underline{V}_n(\theta, \mathbb{P}_n)) - \psi_n(X, \theta, \underline{V}(\theta, P)) \right] \\ (7.16) \quad & + P \left[ \psi_n(X, \theta, \underline{V}(\theta, P)) - \psi(X, \theta, \underline{V}(\theta, P)) \right] = o_p(n^{-1/2}), \end{aligned}$$

where  $\psi_n(X, \theta, \underline{V}(\theta, P))$  equals  $\psi(X, \theta, \underline{V}(\theta, P))$  with the  $\lambda_i$ 's replaced with  $\lambda_{ni}$ 's. To verify that (7.15) is  $o_p(n^{-1/2})$ , we show that  $P[\psi_n(X, \theta, \underline{V}(\theta, P))]$  is Fréchet differentiable. Throughout the proof, the  $\lambda_i$ 's in all expressions are implicitly assumed to be  $\lambda_{ni}$ 's. With  $G^P$  identified with  $P$  and  $G^Q$  identified with  $Q$ , direct calculation shows that the Gateaux derivative  $\partial/\partial\alpha[(1 - \alpha)P + \alpha Q][\psi_n(X, \theta, \underline{V}(\theta, (1 - \alpha)P + \alpha Q))]|_{\alpha=0}$  of  $P[\psi_n(X, \theta, \underline{V}(\theta, P))]$  equals  $(Q -$

$P)\psi_n(X, \theta, \underline{V}(\theta, P))$  plus

$$(7.17) \quad G^P \left[ \sum_{i=1}^s \lambda_{ni} \frac{\dot{w}_i(\theta)}{w_i(\theta)} \tilde{w}_i(\theta) \left[ d_i - \sum_{k=1}^s r_k(\theta) d_k \right] \right] \\ - A(\theta) M^{-}(\theta) G^P \left\{ \left[ d_1 - \sum_{k=1}^s r_k(\theta) d_k \right] \tilde{w}_1(\theta), \dots, \right. \\ \left. \left[ d_s - \sum_{k=1}^s r_k(\theta) d_k \right] \tilde{w}_s(\theta) \right\}^T,$$

where  $d_i \equiv d/d\alpha V_i(\theta, (1 - \alpha)P + \alpha Q)|_{\alpha=0} * V_i^{-1}(\theta, P)$ . The  $d_i$  are found by solving the system of equations

$$\frac{\partial}{\partial \alpha} \left[ \sum_{l=1}^s \lambda_{nl} ((1 - \alpha)G^P + \alpha G^Q) \left[ \left( \frac{w_i(Y, \theta) V_i^{-1}(\theta, (1 - \alpha)G^P + \alpha G^Q)}{\sum_{k=1}^s \lambda_{nk} w_k(Y, \theta) V_k^{-1}(\theta, (1 - \alpha)G^P + \alpha G^Q)} \right) \right. \right. \\ \left. \left. \times \tilde{w}_l(Y, \theta, (1 - \alpha)G^P + \alpha G^Q) \right] \right] \Big|_{\alpha=0} = 0, \\ i = 1, \dots, s - 1,$$

with solution  $d_i = D_{[i, \cdot]}(G^Q - G^P)(\underline{\tilde{w}}(\theta, P))$ , where  $D = \underline{\underline{\lambda}}_n^{-1} M^{-}(\theta, P) M(\theta, P) \underline{\underline{\lambda}}_n$  with  $[i, \cdot]$  denoting the  $i$ th row of the matrix. Substituting the  $d_i$  into (7.17) yields that the Gateaux derivative equals

$$(7.18) \quad (Q - P)\psi_n(X, \theta, \underline{V}(\theta, P)) \\ + G_P \left[ \sum_{i=1}^s \left( \lambda_{ni} \frac{\dot{w}_i(\theta)}{w_i(\theta)} - [A(\theta) M^{-}(\theta)]_{[i, i]} \right) \right. \\ \left. \times \tilde{w}_i(\theta) \left[ D_{[i, \cdot]} - \sum_{k=1}^s r_k(\theta) D_{[k, \cdot]} \right] \right] \\ \times (G^O - G^P)(\underline{\tilde{w}}(\theta, P)).$$

A key ingredient in this calculation is that  $P[-e_I + \underline{r}(\theta)] = 0$ , which implies the complicated term  $d/d\alpha(A(\theta, (1 - \alpha)G^P + \alpha G^Q)M^{-}(\theta, (1 - \alpha)G^P + \alpha G^Q))|_{\alpha=0}$  does not appear in the derivative. Now,

$$(7.19) \quad G_P \left[ \sum_{i=1}^s \left( \lambda_{ni} \frac{\dot{w}_i(\theta)}{w_i(\theta)} - [A(\theta) M^{-}(\theta)]_{[i, i]} \right) \right. \\ \left. \times \tilde{w}_i(\theta) \left[ D_{[i, \cdot]} - \sum_{k=1}^s r_k(\theta) D_{[k, \cdot]} \right] \right] = 0,$$

as can easily be directly verified (by calculating  $A(\theta)$  and  $M^{-}(\theta)$  using (4.7) for the case  $s = 2$ . Thus, the Gateaux derivative equals  $(Q - P)\psi_n(X, \theta, \underline{V}(\theta, P))$ .

To prove (7.15) is  $o_p(n^{-1/2})$ , it is then enough to show that the Gateaux derivative is also a Fréchet derivative with respect to the norm  $\|Q - P\| \equiv \sup\{\sum_{i=1}^s | \int 1_A(y) \tilde{w}_i(y, \theta_0) d(G^Q - G^P)(y) | : A \in \mathbf{B}\}$ . Thus, it suffices to show that

$$(7.20) \quad \left| Q[\psi_n(X, \theta, \underline{V}(\theta, Q))] - P[\psi_n(X, \theta, \underline{V}(\theta, P))] - (Q - P)[\psi_n(X, \theta, \underline{V}(\theta, P))] \right| = o(\|Q - P\|).$$

The expression on the left-hand side of (7.20) equals  $|Q[\psi_n(X, \theta, \underline{V}(\theta, Q)) - \psi_n(X, \theta, \underline{V}(\theta, P))]|$ , and

$$(7.21) \quad \begin{aligned} & \frac{|Q[\psi_n(X, \theta, \underline{V}(\theta, Q)) - \psi_n(X, \theta, \underline{V}(\theta, P))]|}{\|Q - P\|} \\ &= \frac{|Q[\frac{\partial}{\partial \alpha} \psi_n(X, \theta, \underline{V}(\theta, (1 - \alpha)P + \alpha Q))]_{\alpha=0}|}{\|Q - P\|} * O_p(1) \\ &= \left| G_Q \left[ \sum_{i=1}^s \left( \lambda_{ni} \frac{\dot{w}_i(\theta)}{w_i(\theta)} - [A(\theta)M^-(\theta)]_{[i,i]} \right) \right. \right. \\ & \quad \left. \left. \times \tilde{w}_i(\theta) \left[ D_{[i,1]} - \sum_{k=1}^s r_k(\theta) D_{[k,1]} \right] \right] \right| \\ & \quad \times \frac{|(G^Q - G^P)(\tilde{w}(\theta, P))|}{\|Q - P\|} * O_p(1) \\ &= o_p(1) * O_p(1) * O_p(1) = o_p(1) \quad \text{as } \|Q - P\| \rightarrow 0, \end{aligned}$$

using (7.19) to conclude that expression (7.21) is  $o_p(1)$  as  $\|Q - P\| \rightarrow 0$ . Thus (7.15) is  $o_p(n^{-1/2})$ .

To show that (7.16) is  $o_p(n^{-1/2})$ , expand  $\psi_n$  about  $\underline{\lambda} = (\lambda_1, \dots, \lambda_s)^T$ , so that

$$(7.22) \quad \begin{aligned} & P[\psi_n(X, \theta, \underline{V}(\theta, P)) - \psi(X, \theta, \underline{V}(\theta, P))] \\ &= (\underline{\lambda}_n - \underline{\lambda})^T * P \left[ \frac{\partial}{\partial \underline{\lambda}} \psi(X, \theta, \underline{V}(\theta, P)) \Big|_{\underline{\lambda}^*} \right] \end{aligned}$$

for a vector  $\underline{\lambda}^*$  between  $\underline{\lambda}$  and  $\underline{\lambda}_n$ . Since  $\psi$  is the efficient score function for  $\theta$ ,

$$(7.23) \quad \begin{aligned} P \left[ \frac{\partial}{\partial \underline{\lambda}} \psi(X, \theta, \underline{V}(\theta, P)) \Big|_{\underline{\lambda}^*} \right] &= P \left[ \frac{\partial}{\partial \underline{\lambda}} [i_\theta(X) - I_{\theta G} I_G^{-1} i_g(X)] \Big|_{\underline{\lambda}^*} \right] \\ &= P \left[ \frac{\partial}{\partial \underline{\lambda}} i_\theta(X) \Big|_{\underline{\lambda}^*} \right] - \frac{\partial}{\partial \underline{\lambda}} [I_{\theta G} I_G^{-1}] \Big|_{\underline{\lambda}^*} * P \left[ i_g(X) \Big|_{\underline{\lambda}^*} \right] \\ & \quad - [I_{\theta G} I_G^{-1}] \Big|_{\underline{\lambda}^*} * P \left[ \frac{\partial}{\partial \underline{\lambda}} i_g(X) \Big|_{\underline{\lambda}^*} \right]. \end{aligned}$$

Now, if we view the biased sampling density  $p$  of (1.1) as a function of the parameter triplet  $(\underline{\lambda}, \theta, G)$ , then orthogonality of the score for  $\underline{\lambda}$  and the closed linear span of  $(\dot{l}_\theta, \dot{l}_g)$  implies that the first and third terms in the sum (7.23) are zero. The second term in the sum is also zero since score functions have zero expectation. Therefore, (7.16) is  $o_p(n^{-1/2})$ , completing the proof that (5.12) holds.

It remains to verify the stochastic condition of the Master theorem [BKRW, condition (GM1), page 312; equivalent, Van der Vaart (1995), condition (3), Theorem 1]. By Lemma 1 of Van der Vaart (1995), it suffices to show that, for each component  $l = 1, \dots, d$ ,

$$(7.24) \quad \{\psi_{\theta_l} - \psi_{\theta_{0l}}: \|\theta - \theta_0\|_1 < \delta\}$$

is Donsker for some  $\delta > 0$  and

$$(7.25) \quad P_0(\psi_{\theta_l} - \psi_{\theta_{0l}})^2 \rightarrow 0 \quad \text{as } \theta \rightarrow \theta_0.$$

From (5.11), we calculate

$$\begin{aligned} \psi_{\theta_l}(x) - \psi_{\theta_{0l}}(x) &= \left[ \frac{\dot{w}_{il}(y, \theta)}{w_i(y, \theta)} - \frac{\dot{w}_{il}(y, \theta_0)}{w_i(y, \theta_0)} \right] \\ &\quad - \sum_{k=1}^s \left\{ r_k(y, \theta) \left[ \frac{\dot{w}_{kl}(y, \theta)}{w_k(y, \theta)} - \frac{\dot{w}_{kl}(y, \theta_0)}{w_k(y, \theta_0)} \right] \right\} \\ &\quad - \sum_{k=1}^s \left\{ \frac{\dot{w}_{kl}(y, \theta_0)}{w_k(y, \theta_0)} [r_k(y, \theta) - r_k(y, \theta_0)] \right\} \\ &\quad + A_l(\theta) M^-(\theta) \underline{\lambda}^{-1} [-e_i + \underline{r}(y, \theta)] \\ &\quad - A_l(\theta_0) M^-(\theta_0) \underline{\lambda}^{-1} [-e_i + \underline{r}(y, \theta_0)] \\ &\equiv f_{1l\theta}(x) - f_{2l\theta}(x) - f_{3l\theta}(x) + f_{4l\theta}(x) - f_{5l\theta_0}(x), \end{aligned}$$

where the  $1 \times s$  matrix  $A_l(\theta)$  is the  $l$ th row of the matrix  $A(\theta)$  defined in (4.5).

To show that the class  $S_{1l\delta} \equiv \{f_{1l\theta}: \|\theta - \theta_0\|_1 < \delta\}$  is Donsker for some  $\delta > 0$  for each  $l = 1, \dots, d$ , we use the uniform Lipschitz hypothesis (viii) on the weight functions and apply the Jain–Marcus central limit theorem; see, for example, Van der Vaart and Wellner [(1996), pages 213–214]. For fixed  $i \in 1, \dots, s$ ,  $k = 1, \dots, d$ , choose  $\delta$ ,  $K$  and  $\alpha > d/2$  so that (viii) holds. Define  $\Theta_\delta \equiv \{\theta \in \Theta: \|\theta - \theta_0\|_1 < \delta\}$  and  $Z_{nil}(y, \theta) \equiv n^{-1/2}[(\dot{w}_{il}(y, \theta)/w_i(y, \theta)) - (\dot{w}_{il}(y, \theta_0)/w_i(y, \theta_0))]$ , so that  $S_{1l\delta} \subset \cup_{i=1}^s S_{1il\delta}$ , with  $S_{1il\delta} = \{n^{1/2}Z_{nil}(y, \theta): \theta \in \Theta_\delta\}$ . Set  $M_{nk} = n^{-1/2}K$  for  $k = 1, \dots, n$ , so that  $\sum_{k=1}^n EM_{nk}^2 = K^2 = O(1)$ . Define a semimetric  $\rho_\alpha$  on  $\Theta_\delta$  by  $\rho_\alpha(f_{\theta_1}, f_{\theta_2}) = \|\theta_1 - \theta_2\|_1^\alpha$  for  $f_{\theta_1}, f_{\theta_2} \in S_{1il\delta}$ . We require that  $\alpha > d/2$  because this is necessary and sufficient for the semimetric to be Gaussian dominated, that is, to have a finite entropy integral [satisfying the third panel display on page 212 of Van der Vaart and Wellner (1996)]. Next, we verify that the triangular array of norms  $\|Z_{nil}\|_{\Theta_\delta} =$

$\sup_{\theta \in \Theta_\delta} |Z_{nil}(\cdot, \theta)|$  satisfies the Lindeberg condition of Theorem 2.11.11 in Van der Vaart and Wellner (1996):

$$\sum_{i=1}^n E\{\|Z_{nil}\|_{\Theta_\delta}^2\} 1\{\|Z_{nil}\|_{\Theta_\delta} > \eta\} \rightarrow 0 \quad \text{for every } \eta > 0.$$

This follows because (viii) implies  $E\|Z_{nil}\|_{\Theta_\delta}^2 < n^{-1}K^2\delta^{2\alpha} = O(n^{-1})$  and  $1_{[\|Z_{nil}\|_{\Theta_\delta} > \eta]} \leq 1_{[K\delta^\alpha > \eta n^{1/2}]} \rightarrow 0$  as  $n \rightarrow \infty$  for every  $\eta > 0$ .

Since the sequence of covariance functions  $\text{Cov}(Z_{ni_1l}(Y_1, \theta), Z_{ni_2l}(Y_2, \theta))$  converges pointwise in  $\Theta_\delta \times \Theta_\delta$  for each  $i_1 \neq i_2 \in \{1, \dots, s\}$ , the Jain–Marcus CLT yields that, for each fixed  $i = 1, \dots, s, l = 1, \dots, d, S_{1l\delta}$  is Donsker. Thus,  $S_{1l\delta} = \{f_{1l\theta} : \|\theta - \theta_0\|_1 < \delta\}$  is Donsker.

For each  $l = 1, \dots, d$ , the class  $\{f_{2l\theta} : \|\theta - \theta_0\|_1 < \delta\}$  is easily seen to be Donsker, using the assumption that the class  $\mathbf{W}_\delta$  of (5.8) is Donsker and the proof that  $\{f_{1l\theta} : \|\theta - \theta_0\|_1 < \delta\}$  is Donsker.

To show  $S_{4l\delta} \equiv \{f_{4l\theta} : \|\theta - \theta_0\|_1 < \delta\}$  is Donsker, take  $M^-(\theta)$  as in (4.7). Continuous invertibility of  $M_{11}(\theta)$  at  $\theta_0$  and  $\sup_{\|\theta - \theta_0\|_1 < \delta} G(|\dot{w}_{kl}^2(\theta)/w_k(\theta)|^2 \tilde{w}_k(\theta_0)) < \infty$  imply that each element of the vector  $A_l(\theta)M^-(\theta)$  is a bounded function uniformly on  $\Theta_\delta$ . Since each  $\{r_i(\theta) : \|\theta - \theta_0\|_1 < \delta\}$  is Donsker and uniformly bounded, it follows that  $S_{4l\delta}$  is Donsker, as is  $S_{5l\delta} \equiv \{f_{5l\theta_0}\}$ .

Lastly, we show  $S_{3l\delta} \equiv \{f_{3l\theta} : \|\theta - \theta_0\|_1 < \delta\}$  is Donsker. Fix  $k \in \{1, \dots, s\}$ . Define  $d_k(y, \theta) = r_k(y, \theta) - r_k(y, \theta_0)$ . By continuity of  $r_k$  in  $\theta$  at  $\theta_0$ , for any fixed positive  $\epsilon < 1$  we can choose  $\delta$  so that, for each  $y \in \mathbf{Y}$ ,  $\sup_{\|\theta - \theta_0\|_1 < \delta} |d_k(y, \theta)| < \epsilon$ . Write the class  $S_{3l\delta}$  as the sum of the two classes

$$(7.26) \quad \left\{ \frac{\dot{w}_{kl}(y, \theta_0)}{w_k(y, \theta_0)} 1_{[(\dot{w}_{kl}(y, \theta_0)/w_k(y, \theta_0)) \leq \frac{2\epsilon}{1-\epsilon}]} * d_k(y, \theta) : \|\theta - \theta_0\|_1 < \delta \right\},$$

$$(7.27) \quad \left\{ \frac{\dot{w}_{kl}(y, \theta_0)}{w_k(y, \theta_0)} 1_{[(\dot{w}_{kl}(y, \theta_0)/w_k(y, \theta_0)) > \frac{2\epsilon}{1-\epsilon}]} * d_k(y, \theta) : \|\theta - \theta_0\|_1 < \delta \right\}.$$

Since  $\mathbf{W}$  of (5.9) is Donsker, class (7.26) is Donsker as the product of uniformly bounded Donsker classes. Class (7.27) may not be uniformly bounded. We invoke Theorem 2.10.6 of Van der Vaart and Wellner (1996), which applies to Lipschitz transformations of Donsker classes. It asserts that if

$$\begin{aligned} & \left| \frac{\dot{w}_{kl}(y, \theta_0)}{w_k(y, \theta_0)} 1_{[(\dot{w}_{kl}(y, \theta_0)/w_k(y, \theta_0)) > \frac{2\epsilon}{1-\epsilon}]} [d_k(y, \theta_1) - d_k(y, \theta_2)] \right|^2 \\ & \leq \left| \frac{\dot{w}_{kl}(y, \theta_0)}{w_k(y, \theta_0)} 1_{[(\dot{w}_{kl}(y, \theta_0)/w_k(y, \theta_0)) > \frac{2\epsilon}{1-\epsilon}]} - d_k(y, \theta_1) \right|^2 \\ & \quad + \left| \frac{\dot{w}_{kl}(y, \theta_0)}{w_k(y, \theta_0)} 1_{[(\dot{w}_{kl}(y, \theta_0)/w_k(y, \theta_0)) > \frac{2\epsilon}{1-\epsilon}]} - d_k(y, \theta_2) \right|^2 \end{aligned}$$

for all  $\theta_1, \theta_2 \in \Theta_\delta$ , then class (7.27) is Donsker. This follows by straightforward algebra, so that the  $k$ th summand of  $S_{3l\delta}$  is Donsker. Then  $S_{3l\delta}$  itself is Donsker, and condition (7.24) holds.

It remains to verify (7.25). Hypotheses (iv), (vii) and (viii) imply that  $P_0[\sup_{\{\|\theta-\theta_0\|_1<\delta\}}(\psi_{\theta_l}-\psi_{\theta_{0l}})^2] < \infty$  for each  $l = 1, \dots, d$ . Thus, in view of the dominated convergence theorem, it suffices to show that  $\psi_{\theta_l} \rightarrow \psi_{\theta_{0l}}$  pointwise in  $x$ , for each  $l$ . This is true by continuity of the weights and scores for weights in  $\theta$  at  $\theta_0$ .  $\square$

PROOF OF THEOREM 5.4. For  $h \in \mathbf{H}$ , write

$$(7.28) \quad \mathbb{Z}_n(h, \hat{\theta}_n) = \sqrt{n}(\mathbb{G}_n(h, \hat{\theta}_n, \underline{\mathbb{V}}_n(\hat{\theta}_n)) - \mathbb{G}_n(h, \theta_0, \underline{\mathbb{V}}_n(\theta_0)))$$

$$(7.29) \quad + \sqrt{n}(\mathbb{G}_n(h, \theta_0, \underline{\mathbb{V}}_n(\theta_0)) - \mathbb{G}_n(h, \theta_0, \underline{\mathbb{V}}(\theta_0))).$$

Since  $G(h_e^2 r(\theta_0)) < \infty$  and the class  $\mathbf{F}$  is Donsker for each  $F_{i0}$ , piece (7.29) satisfies

$$\|\sqrt{n}(\mathbb{G}_n(h, \theta_0, \underline{\mathbb{V}}_n(\theta_0)) - G(h, \theta_0, \underline{\mathbb{V}}(\theta_0))) - Z(h)\|_{\mathbf{H}} \rightarrow_p 0$$

by Theorem 2.2 of GVW. Expanding, write piece (7.28) as

$$(7.30) \quad \begin{aligned} & \sqrt{n}(\mathbb{G}_n(h, \hat{\theta}_n, \underline{\mathbb{V}}_n(\hat{\theta}_n)) - \mathbb{G}_n(h, \theta_0, \underline{\mathbb{V}}_n(\theta_0))) \\ & = (\nabla \mathbb{G}_n(h, \theta^*, \underline{\mathbb{V}}_n(\theta^*)))^T * \sqrt{n}(\hat{\theta}_n - \theta_0) \end{aligned}$$

for some  $\theta^*$  between  $\theta_0$  and  $\hat{\theta}_n$ , where  $\nabla \mathbb{G}_n(h, \theta, \underline{\mathbb{V}}_n(\theta)) \equiv ((\partial/\partial\theta'_1)\mathbb{G}_n(h, \theta', \underline{\mathbb{V}}_n(\theta'))|_{\theta'=\theta}, \dots, (\partial/\partial\theta'_d)\mathbb{G}_n(h, \theta', \underline{\mathbb{V}}_n(\theta'))|_{\theta'=\theta})^T$ . From representation (1.15) of  $\mathbb{G}_n$ , the gradient  $\nabla \mathbb{G}_n(h, \theta, \underline{\mathbb{V}}_n(\theta))$  can be calculated by the chain rule, using  $(\partial/\partial\theta')\underline{\mathbb{V}}_n(\theta')|_{\theta} = A_n(\theta)M_n^-(\theta)\lambda_n^{-1}\underline{\mathbb{V}}_n(\theta)$ , and equals

$$\begin{aligned} & \frac{-\mathbb{F}_n(h \sum_{i=1}^s r_{ni}(\theta) \frac{\dot{w}_i(Y, \theta)}{w_i(Y, \theta)} r_n(\theta))}{\mathbb{F}_n(r_n(\theta))} + \frac{\mathbb{F}_n(\sum_{i=1}^s r_{ni}(\theta) \frac{\dot{w}_i(Y, \theta)}{w_i(Y, \theta)} r_n(\theta))}{\mathbb{F}_n(r_n(\theta))} * \frac{\mathbb{F}_n(hr_n(\theta))}{\mathbb{F}_n(r_n(\theta))} \\ & + A_n(\theta)M_n^-(\theta) \left[ \frac{-\mathbb{F}_n(h\tilde{w}(\theta)r_n^2(\theta))}{\mathbb{F}_n(r_n(\theta))} + \frac{\mathbb{F}_n(\tilde{w}(\theta)r_n^2(\theta))}{\mathbb{F}_n(r_n(\theta))} * \frac{\mathbb{F}_n(hr_n(\theta))}{\mathbb{F}_n(r_n(\theta))} \right]. \end{aligned}$$

By uniform consistency of  $\mathbb{G}_n(h, \hat{\theta}_n, \underline{\mathbb{V}}_n(\hat{\theta}_n))$ ,  $\nabla \mathbb{G}_n(h, \theta^*, \underline{\mathbb{V}}_n(\theta^*))$  converges in probability to

$$(7.31) \quad \begin{aligned} & -G\left(\sum_{k=1}^s r_k(\theta_0) \frac{\dot{w}_k(\theta_0)}{w_k(\theta_0)} (h - G(h))\right) \\ & - A(\theta_0)M^-(\theta_0)G[r(\theta_0)\tilde{w}(\theta_0)(h - G(h))] \\ & = -G(h\alpha_*(\theta_0)). \end{aligned}$$

Thus for fixed  $h \in \mathbf{H}$ , piece (7.28) converges in distribution to  $-G(h\alpha_*(\theta_0)) * Z_{\theta_0}$ . We then have

$$\begin{aligned} & (\nabla \mathbb{G}_n(h, \theta^*, \underline{\mathbb{V}}_n(\theta^*)))^T * \sqrt{n}(\hat{\theta}_n - \theta_0) \\ & = [(\nabla \mathbb{G}_n(h, \theta^*, \underline{\mathbb{V}}_n(\theta^*)))^T + G(h\alpha_*^T(\theta_0))] * \sqrt{n}(\hat{\theta}_n - \theta_0) \\ & - G(h\alpha_*^T(\theta_0)) * \sqrt{n}(\hat{\theta}_n - \theta_0) \end{aligned}$$

$$\begin{aligned}
&= o_p(1) * O_p(1) - G(h\alpha_\star^T(\theta_0)) * \sqrt{n}(\hat{\theta}_n - \theta_0) \\
&\Rightarrow -G(h\alpha_\star^T(\theta_0)) * Z_{\theta_0} \quad \text{in } l^\infty(\mathbf{H}).
\end{aligned}$$

To establish the convergence of  $Z'_n(h, \hat{\theta}_n)$  to  $Z'(h)$  of (5.6), it remains to show that  $\lim_n Z'_n$  has covariance (5.7). Using (7.30) and (7.31), write

$$\begin{aligned}
Z'_n(h, \hat{\theta}_n) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \begin{pmatrix} \tilde{l}_{\theta_0}(X_i) \\ -G(h\alpha_\star^T(\theta_0)) * \tilde{l}_{\theta_0}(X_i) + \tilde{l}_{G^0}(\pi_h)(X_i) \end{pmatrix} + o_p(1) \\
&\rightarrow_d N_{d+1} \left( \mathbf{0}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12}(h) \\ \Sigma_{21}(h) & \Sigma_{22}(h, h) \end{bmatrix} \right),
\end{aligned}$$

where  $\tilde{l}_\theta$  is the efficient influence function for  $\theta$  and  $\tilde{l}_{G^0}$  is the efficient influence function for  $G$  when  $\theta = \theta_0$  is known. Evidently,  $\Sigma_{11} = I^{-1}(\theta_0)$ . Since the scores for  $\theta$  must be orthogonal to the nuisance parameter scores, and  $\tilde{l}_{G^0}(\pi_h)$  is in the space of nuisance parameter scores, it follows that  $E(\tilde{l}_{\theta_0}(X)\tilde{l}_{G^0}(\pi_h)(X)) = \mathbf{0}$ . Then elementary calculation yields  $\Sigma_{12}(h) = -I^{-1}(\theta_0)G(h\alpha_\star(\theta_0))$  and  $\Sigma_{22}(h, h) = G(h\alpha_\star^T(\theta_0))I^{-1}(\theta_0)G(h\alpha_\star(\theta_0)) + I_{G^0}^{-1}(h, h) = I_G^{-1}(h, h)$ . It follows that  $\text{AsymCov}(Z'_n(h_1, \hat{\theta}_n), Z'_n(h_2, \hat{\theta}_n))$ , for  $h_1, h_2 \in \mathbf{H}$ , matches the covariance formula (5.7).  $\square$

**PROOF OF PROPOSITION 5.2.** Adding and subtracting  $\mathbb{G}_n(\underline{w}(\hat{\theta}_n), \theta_0)$  and  $\mathbb{G}_n(\underline{w}(\theta_0), \theta_0)$ , and expanding about  $\theta_0$ , allows us to write

$$\begin{aligned}
(7.32) \quad \sqrt{n}(\underline{W}_n - \underline{W}_0) &= [\nabla \mathbb{G}_n(\underline{w}(\hat{\theta}_n), \theta^*, \underline{V}_n(\theta^*)) + \mathbb{G}_n(\underline{w}(\theta^*), \theta_0, \underline{V}_n(\theta_0))]^T \\
&\quad \times \sqrt{n}(\hat{\theta}_n - \theta_0)
\end{aligned}$$

$$(7.33) \quad + \sqrt{n}(\underline{W}_n(\theta_0) - \underline{W}_0)$$

for some  $\theta^*$  between  $\theta_0$  and  $\hat{\theta}_n$ . By Proposition 2.2 of GVW, the process in (7.33) converges in distribution to the first two lines of (5.13), with covariance given by the first three lines of (5.14). Slutsky's theorem and (7.31) show that expression (7.32) converges to the third and fourth lines of (5.13), with covariance given by the third and fourth lines of (5.14). Since  $\underline{W}_n(\theta_0)$  has influence function  $\tilde{l}_{\underline{W}_0} = [\tilde{l}_{G^0}(\pi_{w_1(\theta_0)}), \dots, \tilde{l}_{G^0}(\pi_{w_s(\theta_0)})]^T$ , the asymptotic covariance of the two terms (7.32) and (7.33) equals

$$\begin{aligned}
&\underline{W}_0 \left[ -G(\underline{w}(\theta_0)\alpha_\star^T(\theta_0)) + G \left( \frac{\dot{w}_1(\theta_0)}{w_1(\theta_0)} \tilde{w}_1(\theta_0), \dots, \frac{\dot{w}_s(\theta_0)}{w_s(\theta_0)} \tilde{w}_s(\theta_0) \right)^T \right] \\
&\quad \times E[\tilde{l}_{\theta_0}(X)\tilde{l}_{\underline{W}_0}^T(X)].
\end{aligned}$$

Since  $E(\tilde{l}_{\theta_0}(X)\tilde{l}_{G^0}(\pi_h)) = \mathbf{0}$  for  $h \in \mathbf{H}$ , and each  $w_i(\cdot, \theta_0) \in \mathbf{H}$  by requirement (5.8), the last expectation term is zero.  $\square$



PROOF OF PROPOSITION 5.3. Write

$$(7.34) \quad \sqrt{n}(\hat{F}_{ni}(h, \hat{\theta}_n) - F_{i0}(h)) = \sqrt{n} \left( \frac{\mathbb{G}_n(hw_i(\hat{\theta}_n), \hat{\theta}_n, \underline{V}_n(\hat{\theta}_n))}{\mathbb{G}_n(w_i(\hat{\theta}_n), \hat{\theta}_n, \underline{V}_n(\hat{\theta}_n))} - \frac{\mathbb{G}_n(hw_i(\theta_0), \theta_0, \underline{V}_n(\theta_0))}{\mathbb{G}_n(w_i(\theta_0), \theta_0, \underline{V}_n(\theta_0))} \right) \\ (7.35) \quad + \sqrt{n}(\hat{F}_{ni}(h, \theta_0) - F_{i0}(h)).$$

The process of (7.35) converges in distribution to the first two lines of (5.15) by Theorem 2.3 of GVW. Using the identity  $(\partial/\partial\theta')\tilde{w}_i(y, \theta')|_{\theta'=\theta} = \dot{l}_\theta(i, y)$ , a Taylor expansion shows that piece (7.34) equals

$$[\nabla\mathbb{G}_n(h\tilde{w}_i(\hat{\theta}_n, \mathbb{F}_n), \theta^*, \underline{V}_n(\theta^*)) + \mathbb{G}_n(h\dot{l}_{\theta^*}(i, Y), \theta_0, \underline{V}_n(\theta_0))]^T * \sqrt{n}(\hat{\theta}_n - \theta_0) + o_p(1)$$

for some  $\theta^*$  between  $\theta_0$  and  $\hat{\theta}_n$ , which converges in distribution to the term in the last line of (5.15). From the expansion used by GVW [page 1107], the influence function  $\tilde{l}_{F_{i0}}$  of  $\hat{F}_{ni}$  equals

$$\tilde{l}_{F_{i0}}(\pi_h) = \tilde{l}_{G^0}(\pi_{w_i(\theta_0)h}) + G(\tilde{w}_i(\theta_0)h)W_{0i}^{-1}\tilde{l}_{G^0}(\pi_{w_i(\theta_0)}).$$

The asymptotic covariance between (7.34) and (7.35) equals

$$[-G([\tilde{w}_i(\theta_0) - G(\tilde{w}_i(\theta_0))]a_*(\theta_0)) + G(h\dot{l}_{\theta_0}(i, Y))]^T E[\tilde{l}_{\theta_0}(X)\tilde{l}_{F_{i0}(\theta_0)}(\pi_h)],$$

which is zero by the assumption that each  $w_i(\cdot, \theta_0)$  and each  $w_i(\cdot, \theta_0)h$  is in  $\mathbf{H}$ .  $\square$

PROOF OF THEOREM 6.1, COROLLARY 6.1 AND COROLLARY 6.2. Since the weights are positive, the connectivity of  $\mathbf{G}^*(\theta)$  for all  $\theta$  in a neighborhood of  $\theta_0$  holds automatically. As discussed in Section 3 of Gilbert, Lele and Vardi (1999), the remaining conditions of identifiability Theorem 1.1 are met if  $h_{ik}(0) = 0$  for each  $i = 1, \dots, s - 1, k = 1, \dots, d$ , and for some  $i \in \{1, \dots, s - 1\}$ , the set of functions  $\{h_{i1}, \dots, h_{id}\}$  is linearly independent. The MLE obtained by procedure (1)–(3) exists uniquely in the limit with probability 1, as verified in Example 6.1. In addition, methods for verifying the hypotheses of the theorems in Section 5 were outlined in Example 6.1.  $\square$

**8. Discussion and further problems.** Results from a simulation study of the MLE  $(\hat{\theta}_n, \mathbb{G}_n)$  in the two- and three-sample GLR models [Gilbert (1996); Gilbert, Lele and Vardi (1999)] corroborate the large sample properties described here. The simulation design utilized a quasi-Newton scheme to compute the MLE via procedure (1)–(3). The log profile partial likelihood was found to be smooth and strictly concave for the great majority of data configurations. Asymptotic unbiasedness and normality of  $(\hat{\theta}_n, \mathbb{G}_n)$  were confirmed, as were consistency of the plug-in and bootstrap variance estimates. Further, the likelihood ratio, Wald and efficient score tests were found to be consistent and possess approximate chi-squared distributions, and profile partial likelihood-based confidence sets for  $\theta_0$  had accurate coverage probabilities.

In summary, the maximum likelihood estimator in the semiparametric biased sampling model shares many properties with the maximum likelihood estimator in Cox's semiparametric proportional hazards model. The estimate  $\hat{\theta}_n$  is computed conveniently by maximizing a smooth quadratic log profile partial likelihood, and then the NPMLE of the baseline distribution is computed. The asymptotic properties of the joint MLEs are comparable, each converging at rate  $n^{1/2}$  and achieving semiparametric efficiency, with similar forms and derivations of information bounds. These similarities are not surprising, as the  $s$ -group proportional hazards model, defined by  $\lambda(y, \theta|i) = \exp(\theta_i)\lambda(y|s)$  for  $i = 1, \dots, s$ ,  $\theta = (\theta_1, \dots, \theta_s)^T$ ,  $\theta_s \equiv 0$ , has the analytic form of an  $s$ -sample biased sampling model (albeit with weight functions depending on the infinite-dimensional parameter  $G$ ) with  $w_i(y, \theta, G) = (1 - G(y))^{\exp(\theta_i)-1}$ , where  $G$  is the distribution function  $Y$  for the  $s$ th group.

**PROBLEM 8.1** (All weight functions unknown). As alluded to in the Introduction, the one-sample biased sampling model (with no restrictions on  $G$ ) is rarely identifiable when the weight function depends on  $\theta$ . In cases that it is [see Gilbert, Lele and Vardi (1999), Theorem 1], no systematic estimation procedure exists. In multiple-sample biased sampling models in which all weight functions depend on an unknown parameter, the class of identifiable models is larger [see Gilbert, Lele and Vardi (1999), Theorem 3]. If one of the normalizing constants is known, the methodology given here applies. If not, which is commonly the case, there does not exist a systematic procedure for estimation. A thorough treatment of these situations would be of interest.

**PROBLEM 8.2** (Adaptive or nonparametric estimation of the weight functions). In many problems, there may be little rationale for choosing a parametric form for the weight functions. For instance, consider the vaccine trial Examples 1.1, 1.2, and 6.1. For vaccines against some heterogeneous pathogens, including HIV, there are very few data available to guide selection of the weight functions. In this situation, it is of interest to consider more adaptive estimation of the  $\{w_i\}$ . In the econometrics literature, Manski (1993) and references therein consider estimation of the weight function in one biased sample via kernel methods. Sun and Woodroffe (1997) consider one- and two-sample biased sampling models in which the underlying distribution function is specified parametrically,  $w_2 = 1$  if there is a second sample, and the weight function  $w_1$  is estimated nonparametrically. They assume  $w_1$  is a monotone density function and construct a consistent (penalized) maximum likelihood estimator  $\hat{w}_1$ . It would be of interest to investigate the asymptotic distributional properties of this estimator. Another approach is to estimate  $w_1$  by smoothing splines through maximization of a roughness-penalized version of the profile partial likelihood considered here. We are pursuing research in this direction.

Adaptive or nonparametric estimation of the weight functions is of interest in its own right, or as a goodness-of-fit diagnostic for checking the parametric form of the weight functions.

**Acknowledgment.** I am grateful to Jon Wellner for his very helpful and kind guidance.

## REFERENCES

- AGRESTI, A. (1984). *Analysis of Ordinal Categorical Data*. Wiley, New York.
- ANDERSEN, P. K. and GILL, R. D. (1982). Cox's regression model for counting processes: a large sample study. *Ann. Statist.* **10** 1100–1120.
- BEGUN, J. M., HALL, W. J., HUANG, W-M. and WELLNER, J. A. (1983). Information and asymptotic efficiency in parametric-nonparametric models. *Ann. Statist.* **11** 432–452.
- BICKEL, P. J., KLAASSEN, C. A., RITOV, Y. and WELLNER, J. A. (1993). *Efficient and Adaptive Estimation for Semiparametric Models*. Johns Hopkins Univ. Press, Baltimore, MD.
- COSSLETT, S. R. (1981). Maximum likelihood estimator for choice based samples. *Econometrika* **49** 1289–1316.
- COX, D. R. and SNELL, E. J. (1989). *The Analysis of Binary Data*. 2nd ed. Chapman and Hall, London.
- DUDLEY, R. M. (1984). A course on empirical process. *Ecole d' Été de Probabilités de Saint Flour XII. Lecture Notes in Math.* **1097** 1–142. Springer, New York.
- DUDLEY, R. M. (1985). An extended Wichura theorem, definitions of Donsker classes, and weighted empirical processes. *Probability in Banach Spaces V. Lecture Notes in Math.* **1153** 1306–1326. Springer, New York.
- FIENBERG, S. E. (1980). *The Analysis of Cross-Classified Categorical Data*. MIT Press.
- GILBERT, P. B. (1996). Sieve analysis: statistical methods for assessing differential vaccine protection against HIV types. Ph.D. dissertation, Univ. Washington.
- GILBERT, P. B., LELE, S. R. and VARDI, Y. (1999). Maximum likelihood estimation in semiparametric selection bias models with application to AIDS vaccine trials. *Biometrika* **86** 27–43.
- GILBERT, P. B., SELF, S. G. and ASHBY, M. A. (1998). Statistical methods for assessing differential vaccine protection against HIV types. *Biometrics* **54** 799–814.
- GILL, R. D., VARDI, Y. and WELLNER, J. A. (1988). Large sample theory of empirical distributions in biased sampling models. *Ann. Statist.* **16** 1069–1112.
- HUBER, P. J. (1967). The behavior of maximum likelihood estimates under nonstandard conditions. *Proceedings of the Fifth Berkeley Symp. Math. Statist. Probab.* **1** 221–233. Univ. California Press, Berkeley.
- KIEFER, J. and WOLFOWITZ, J. (1956). Consistency of the maximum likelihood estimator in the presence of infinitely many incidental parameters. *Ann. Math. Statist.* **27** 887–906.
- MANSKI, C. F. (1993). The selection problem in econometrics and statistics. In *Handbook of Statistics* **11** (G. S. Maddala, C. R. Rao and H. D. Vinod, Eds.) 73–84. North-Holland, Amsterdam.
- MANSKI, C. F. and LERMAN, S. R. (1977). The estimation of choice probabilities from choice-based samples. *Econometrics* **45** 1977–1988.
- OSSIANDER, M. (1987). A central limit theorem under metric entropy with  $L_2$  bracketing. *Ann. Probab.* **15** 897–919.
- POLLARD, D. (1984). *Convergence of Stochastic Processes*. Springer, New York.
- QIN, J. (1998). Inferences for case-control and semiparametric two-sample density ratio models. *Biometrika* **85** 619–630.
- SUN, J. and WOODROOFE, M. (1997). Semiparametric estimates for biased sampling models. *Statist. Sinica* **7** 545–575.
- TRICOMI, F. G. (1957). *Integer Equations*. Interscience, New York.
- VAN DER VAART, A. (1994). Bracketing smooth functions. *Stochast. Process. Appl.* **52** 93–105.
- VAN DER VAART, A. (1995). Efficiency of infinite dimensional  $M$ -estimators. *Statist. Neerlandica* **49** 9–30.
- VAN DER VAART, A. and WELLNER, J. A. (1996). *Weak Convergence and Empirical Processes*. Springer, New York.

- VARDI, Y. (1985). Empirical distributions in selection bias models. *Ann. Statist.* **13** 178–203.
- WELLNER, J. A. and ZHAN, Y. (1998). Bootstrapping  $Z$ -estimators. Technical Report 308, Dept. Statistics, Univ. Washington.

DEPARTMENT OF BIostatISTICS  
HARVARD SCHOOL OF PUBLIC HEALTH  
BOSTON, MASSACHUSETTS 02115  
E-MAIL: [pgilbert@hsph.harvard.edu](mailto:pgilbert@hsph.harvard.edu)