

# Empirical likelihood approaches under complex sampling designs

Y.G. Berger

Southampton Statistical Sciences Research Institute,  
University of Southampton

This is the submitted version which will be published in final form in Wiley Statsref

<http://wileyactual.com/wileystatsref/>

## Abstract

There are two different empirical likelihood approaches for complex sampling designs: the “*pseudoempirical likelihood*” introduced by Chen and Sitter (1999) and “*unequal probability empirical likelihood*” approach proposed by Berger and Torres (2016). Both approaches are described and reviewed critically. The key difference is the fact that the self-normalisation property of the pseudoempirical likelihood approach is limited to unidimensional parameters. This property holds for multidimensional parameters, with the unequal probability empirical likelihood approach. This is a brief description of the key empirical likelihood approaches for complex sampling. This is not an exhaustive account of all the applications of empirical likelihood in survey sampling.

*Keywords:* Design-based approach, estimating equations, inclusion probabilities, side information, stratification

# 1 Introduction

Consider a finite population  $\mathcal{U} = \{1, \dots, N\}$  of  $N$  units. For each unit  $i \in \mathcal{U}$ , we have a vector  $\mathbf{y}_i$  of variables, where  $\mathbf{y}_i \in \mathbb{R}^{d_y}$ . We consider a “*design-based approach*”; that is, we assume that  $\mathbf{y}_i$  are vectors of constants (Neyman, 1938). A random sample  $\mathcal{S}$  is selected from  $\mathcal{U}$  according to a sampling design, which specifies the probability distribution of  $\mathcal{S}$ . The sampling distribution is therefore given by the sampling design. The design-based approach is the core of survey sampling theory. It offers a robust non-parametric approach for survey data, which does not rely on distributional assumption about  $\mathbf{y}_i$ .

We consider that the population  $\mathcal{U}$  is split into  $H$  non-overlapping strata  $\mathcal{U}_1, \dots, \mathcal{U}_h, \dots, \mathcal{U}_H$  such that  $\cup_{h=1}^H \mathcal{U}_h = \mathcal{U}$ . Within each stratum  $\mathcal{U}_h$ , a sample of  $n_h$  units is selected with unequal selection probabilities  $\pi_i$ . The overall sample size is the constant  $n = \sum_{h=1}^H n_h$ . With and without replacement sampling will be considered. We shall focus on single-stage designs.

Empirical likelihood has its origin in survey sampling theory. One of the first attempts to formulate a likelihood-based approach in survey sampling is due to Godambe (1966) who showed that under the design-based approach, the likelihood function is flat and cannot be used for inference. In order to solve this problem, Hartley and Rao (1968) developed the first application of an empirical likelihood approach under equal probability sampling. Owen (1988) popularised this approach into the mainstream statistics (see also Owen, 2001). There have been many recent developments of empirical likelihood based methods in survey sampling. Most of which can be classified into two categories: the unequal probability empirical likelihood approach (Berger and Torres, 2016) and the pseudoempirical likelihood approach (Chen and Sitter, 1999; Wu and Rao, 2006).

# 2 Parameter and side information

Let  $\boldsymbol{\theta}_{\mathcal{U}} \in \mathbb{R}^{d_{\theta}}$  be an unknown vector of  $d_{\theta}$  parameters. For example,  $\boldsymbol{\theta}_{\mathcal{U}}$  can be the parameter of a generalised linear regression model. Let  $\boldsymbol{\varphi}_{\mathcal{U}} \in \mathbb{R}^{d_{\varphi}}$  denote another vector called “*side information*” which is assumed known without sampling errors. For example,  $\boldsymbol{\varphi}_{\mathcal{U}}$  can be population-level means, counts or proportions, from large external censuses or surveys (see Example 2.1). Side information is often used in survey sampling (e.g. Kott, 2009), but not limited to this field. It can also be found in the empirical likelihood literature (Owen, 1991, 2001, §3.10) and in the econometric literature (Imbens and Lancaster, 1994). Taking into account of side information may increase the efficiency (Owen, 1991).

Let  $\boldsymbol{\psi}_{\mathcal{U}} = (\boldsymbol{\theta}_{\mathcal{U}}^{\top}, \boldsymbol{\varphi}_{\mathcal{U}}^{\top})^{\top}$ . We assume that  $\boldsymbol{\psi}_{\mathcal{U}}$  is the unique solution to the population multidimensional estimating equation (Godambe, 1960)

$$\sum_{i \in \mathcal{U}} \mathbf{g}(\mathbf{y}_i, \boldsymbol{\theta}, \boldsymbol{\varphi}) = \mathbf{0}_{d_g}, \quad \text{if and only if } \boldsymbol{\theta} = \boldsymbol{\theta}_{\mathcal{U}} \text{ and } \boldsymbol{\varphi} = \boldsymbol{\varphi}_{\mathcal{U}}; \quad (1)$$

where

$$\mathbf{g}(\mathbf{y}_i, \boldsymbol{\theta}, \boldsymbol{\varphi}) := \{\mathbf{e}(\mathbf{y}_i, \boldsymbol{\theta}, \boldsymbol{\varphi})^{\top}, \mathbf{f}(\mathbf{y}_i, \boldsymbol{\varphi})^{\top}\}^{\top} \in \mathbb{R}^{d_g} \quad (d_g \geq d_{\theta}), \quad (2)$$

with  $\mathbf{f}(\mathbf{y}_i, \boldsymbol{\varphi}_U) \in \mathbb{R}^{d_f}$  ( $d_f \geq d_\varphi$ ) and  $\mathbf{e}(\mathbf{y}_i, \boldsymbol{\theta}, \boldsymbol{\varphi}) \in \mathbb{R}^{d_e}$  ( $d_e \geq d_\theta$ ). Here,  $\mathbf{0}_{d_g}$  denotes a  $d_g$ -vector of zeros. The estimating function  $\mathbf{e}(\mathbf{y}_i, \boldsymbol{\theta}, \boldsymbol{\varphi})$  defines  $\boldsymbol{\theta}_U$ . For example,  $\boldsymbol{\theta}_U$  can be the parameter of a generalised linear regression model (Chen and Van Keilegom, 2009). Under a design-based approach, the parameter  $\boldsymbol{\theta}_U$  is a vector of unknown constants. Note that the function  $\mathbf{e}(\mathbf{y}_i, \boldsymbol{\theta}, \boldsymbol{\varphi})$  may or may not depend on  $\boldsymbol{\varphi}$ .

Note that (1) implies that  $\boldsymbol{\varphi}_U$  is such that

$$\sum_{i \in \mathcal{U}} \mathbf{f}(\mathbf{y}_i, \boldsymbol{\varphi}_U) = \mathbf{0}, \quad (3)$$

by definition. The  $\mathbf{f}(\mathbf{y}_i, \boldsymbol{\varphi}_U)$  are usually a function of a sub-vector  $\tilde{\mathbf{y}}_i$  of  $\mathbf{y}_i$ , with  $\tilde{\mathbf{y}}_i$  known  $\forall i \in \mathcal{U}$ . We assume that  $\mathbf{f}(\mathbf{y}_i, \boldsymbol{\varphi}_U)$  are known  $\forall i \in \mathcal{S}$ . The  $\mathbf{f}(\mathbf{y}_i, \boldsymbol{\varphi}_U)$  are called ‘auxiliary variables’ in the survey sampling literature (e.g. Deville and Särndal, 1992). We shall treat  $\boldsymbol{\varphi}_U$  as a vector of constants, not as a parameter to estimate.

The estimating function  $\mathbf{f}(\mathbf{y}_i, \boldsymbol{\varphi}_U)$  cannot be a function of  $\boldsymbol{\theta}_U$ , because  $\boldsymbol{\varphi}_U$  is assumed known and defined by (3). Indeed, if  $\mathbf{f}(\mathbf{y}_i, \boldsymbol{\varphi})$  depends on  $\boldsymbol{\theta}_U$ , the vector  $\boldsymbol{\varphi}_U$  cannot be known, unless  $\boldsymbol{\theta}_U$  is known, which is a trivial situation where inference is not necessary.

**Example 2.1.** Let  $\mathbf{y}_i = (\mathbf{x}_i^\top, \delta_i)^\top$ , where  $\mathbf{x}_i$  is some covariates and  $\delta_i$  is a binary variable:  $\delta_i = 1$  for a success and  $\delta_i = 0$  for a failure. Suppose that we know the population success rate  $\boldsymbol{\varphi}_U = N^{-1} \sum_{i \in \mathcal{U}} \delta_i$ , and we wish to fit a binary logistic model with  $\delta_i$  as dependent variable and  $\mathbf{x}_i$  as covariates. The estimating functions are

$$\begin{aligned} \mathbf{e}(\mathbf{y}_i, \boldsymbol{\theta}, \boldsymbol{\varphi}) &= \mathbf{x}_i^\top \delta_i - \mathbf{x}_i^\top \exp(\mathbf{x}_i^\top \boldsymbol{\theta}) \{1 + \exp(\mathbf{x}_i^\top \boldsymbol{\theta})\}^{-1}, \\ \mathbf{f}(\mathbf{y}_i, \boldsymbol{\varphi}) &= \delta_i - \boldsymbol{\varphi}. \end{aligned}$$

More examples can be found in Imbens and Lancaster (1994), Berger and Torres (2016) and Oğuz-Alper and Berger (2016).

### 3 Unequal probability empirical likelihood

Berger and Torres’s (2016) “empirical log-likelihood function” is defined by

$$\ell_{\max}(\boldsymbol{\theta}, \boldsymbol{\varphi}_U) := \max_{\mathbf{p}} \left\{ \ell(\mathbf{p}) : p_i > 0, \sum_{i \in \mathcal{S}} \frac{p_i}{\pi_i} \mathbf{g}(\mathbf{y}_i, \boldsymbol{\theta}, \boldsymbol{\varphi}_U) = \mathbf{0}_{d_g}, \sum_{i \in \mathcal{S}} p_i \mathbf{z}_i = \frac{\mathbf{n}_H}{n} \right\}, \quad (4)$$

where  $\mathbf{p} = (p_i : i \in \mathcal{S})^\top$  is the  $n$ -vector of  $p_i$ ,

$$\begin{aligned} \ell(\mathbf{p}) &:= \sum_{i \in \mathcal{S}} \log(p_i), \\ \mathbf{z}_i &:= (z_{i1}, \dots, z_{ih}, \dots, z_{iH})^\top, \\ z_{ih} &:= \begin{cases} 1 & \text{if } i \in \mathcal{U}_h, \\ 0 & \text{otherwise,} \end{cases} \\ \mathbf{n}_H &:= \sum_{i \in \mathcal{S}} \mathbf{z}_i = (n_1, \dots, n_h, \dots, n_H)^\top. \end{aligned}$$

The  $\mathbf{z}_i$  are stratification variables and  $\mathbf{n}_H$  is the strata allocation, and  $\ell(\mathbf{p})$  is Owen's (1988) empirical log-likelihood function. Berger and Torres's (2016) empirical log-likelihood function is expressed in term of  $m_i := np_i\pi_i^{-1}$ . If we substitute  $p_i$  within (4) by  $n^{-1}m_i\pi_i$ , straightforward algebra shows (4) reduces to Berger and Torres's (2016) empirical log-likelihood function plus a quantity that does not depend on  $\boldsymbol{\theta}$ . We prefer expressing (4) in term of  $p_i$  in order to ease the comparison with other empirical likelihood approaches.

The constraints within (4) differs from Owen's (1988) constraints. We have  $\pi_i$  within the first constraint involving the parameter. We also have a second stratification constraint not motivated by moment conditions. This constraint implies  $\sum_{i \in \mathcal{S}} p_i = 1$ , since  $\mathbf{1}_H^\top \mathbf{z}_i = 1$  and  $\mathbf{1}_H^\top \mathbf{n}_H = n$ , where  $\mathbf{1}_H$  is the  $H$ -vector of 1. Hence, the stratification constraint can be viewed as a generalisation of Owen's (1988) leading constraint  $\sum_{i \in \mathcal{S}} p_i = 1$ .

The "maximum empirical likelihood estimator"  $\hat{\boldsymbol{\theta}}$  maximises  $\ell_{\max}(\boldsymbol{\theta}, \boldsymbol{\varphi}_U)$  by definition. Berger and Torres (2016) show that  $\hat{\boldsymbol{\theta}}$  is also the solution of the sample estimating equation

$$\sum_{i \in \mathcal{S}} \hat{m}_i(\boldsymbol{\varphi}_U) \mathbf{e}(\mathbf{y}_i, \boldsymbol{\theta}, \boldsymbol{\varphi}_U) = \mathbf{0}_{d_g}, \quad (5)$$

where

$$\begin{aligned} \hat{m}_i(\boldsymbol{\varphi}_U) &:= n \hat{p}_i(\boldsymbol{\varphi}_U) \pi_i^{-1}, \\ \hat{p}_i(\boldsymbol{\varphi}_U) &:= n^{-1} \{1 + \boldsymbol{\eta}(\boldsymbol{\varphi}_U)^\top \mathbf{c}_i(\boldsymbol{\varphi}_U) \pi_i^{-1}\}^{-1}, \\ \mathbf{c}_i(\boldsymbol{\varphi}_U) &:= \{\mathbf{f}(\mathbf{y}_i, \boldsymbol{\varphi}_U)^\top, \pi_i \mathbf{z}_i^\top\}^\top. \end{aligned} \quad (6)$$

Here,  $\boldsymbol{\eta}(\boldsymbol{\varphi}_U)$  is a Lagrangian parameter which is such that

$$\sum_{i \in \mathcal{S}} \hat{m}_i(\boldsymbol{\varphi}_U) \mathbf{c}_i(\boldsymbol{\varphi}_U) = (\mathbf{0}_{d_f}^\top, \mathbf{n}_H^\top)^\top.$$

A modified Newton-Raphson algorithm (e.g. Polyak, 1987) can be used to compute  $\boldsymbol{\eta}(\boldsymbol{\varphi}_U)$ . We assume that  $(\mathbf{0}_{d_f}^\top, \mathbf{n}_H^\top)^\top$  is an inner point of the convex conical hull of  $\{\mathbf{c}_i(\boldsymbol{\varphi}_U) : i \in \mathcal{S}\}$ , so that a unique solution  $\boldsymbol{\eta}(\boldsymbol{\varphi}_U)$  exists.

The quantities  $\hat{m}_i(\boldsymbol{\varphi}_U)$  are the empirical likelihood weights. If we do not have side information,  $\mathbf{f}(\mathbf{y}_i, \boldsymbol{\varphi})$  is removed from (2) and  $\hat{m}_i(\boldsymbol{\varphi}_U)$  reduces to the standard sampling weights  $\pi_i^{-1}$ . If  $\boldsymbol{\theta}_U$  is a population mean, we obtain the Horvitz and Thompson's (1952) estimator with  $\mathbf{e}(\mathbf{y}_i, \boldsymbol{\theta}, \boldsymbol{\varphi}_U) = \mathbf{y}_i - \boldsymbol{\theta} N n^{-1} \pi_i$ . Note that the calibration property (Deville and Särndal, 1992) holds because of  $\sum_{i \in \mathcal{S}} \hat{m}_i(\boldsymbol{\varphi}_U) \mathbf{f}(\mathbf{y}_i, \boldsymbol{\varphi}_U) = \mathbf{0}_{d_f}$  and (3). We have this property because we maximise (4) and  $\boldsymbol{\varphi}_U$  is constant. In survey sampling literature, calibration is viewed as a weighting procedure, rather than the consequence of the maximisation of likelihood function.

On the main advantage of this §'s approach is that the function (4) can be used for testing. Suppose we wish to test  $H_0 : \boldsymbol{\theta}_U^{(1)} = \boldsymbol{\theta}_0^{(1)}$ , where  $\boldsymbol{\theta}_U^{(1)} \in \mathbb{R}^{d_{\boldsymbol{\theta}^{(1)}}}$  is a sub-parameter of  $\boldsymbol{\theta}_U$ ; that is,  $\boldsymbol{\theta}_U = (\boldsymbol{\theta}_U^{(1)\top}, \boldsymbol{\theta}_U^{(2)\top})^\top$ . Oğuz-Alper and Berger (2016) showed that under  $H_0$ ,

$$\hat{r}(\boldsymbol{\theta}^{(1)}, \boldsymbol{\varphi}_U) := 2 \left\{ \ell_{\max}(\hat{\boldsymbol{\theta}}, \boldsymbol{\varphi}_U) - \max_{\boldsymbol{\theta}_U^{(2)} \in \Theta^{(2)}} \ell_{\max}(\boldsymbol{\theta}, \boldsymbol{\varphi}_U) \right\} \xrightarrow{d} \chi_{d_{\boldsymbol{\theta}^{(1)}}}^2, \quad \text{if } \boldsymbol{\theta}^{(1)} = \boldsymbol{\theta}_0^{(1)}, \quad (7)$$

under with replacement stratified sampling, as  $n \rightarrow \infty$ . Inverse testing can be used to construct confidence intervals, when  $\boldsymbol{\theta}_U^{(1)}$  is unidimensional. The property (7) is also known as the self-normalisation property. Oğuz-Alper and Berger (2016) also showed how this approach can be extended to multi-stage designs.

### 3.1 Extension for large sampling fractions

It is common practice for business surveys to use non-negligible sampling fractions  $n/N$ . In this case, sampling without-replacement is preferred over sampling with-replacement, in order to avoid selecting the same units several times. With large sampling fractions and without-replacement sampling, the maximum empirical likelihood estimator is still the solution of (5), but the property (7) does not hold. Berger and Torres (2016) proposed a “*penalised empirical likelihood function*” as a solution to this problem. This function is defined by

$$\begin{aligned} \tilde{\ell}_{\max}(\boldsymbol{\theta}, \boldsymbol{\varphi}_U) &:= \max_{\mathbf{p}} \left\{ \tilde{\ell}(\mathbf{p}) : p_i > 0, \sum_{i \in \mathcal{S}} (p_i q_i - \psi_i) \frac{1}{\pi_i} \mathbf{g}(\mathbf{y}_i, \boldsymbol{\theta}, \boldsymbol{\varphi}_U) = \mathbf{0}_{d_g}, \right. \\ &\quad \left. \sum_{i \in \mathcal{S}} (p_i q_i - \psi_i) \mathbf{z}_i = \frac{\mathbf{n}_H}{n} \right\}, \end{aligned} \quad (8)$$

where

$$\begin{aligned} \tilde{\ell}(\mathbf{p}) &:= \sum_{i \in \mathcal{S}} \log(p_i) - n \sum_{i \in \mathcal{S}} p_i + n, \\ q_i &:= (1 - \pi_i)^{1/2}, \\ \psi_i &:= (q_i - 1)n^{-1}. \end{aligned}$$

The penalty  $n \sum_{i \in \mathcal{S}} p_i + n$  is necessary for the pivotal property (9) to hold. The  $q_i$  are Hájek’s (1964) finite population corrections.

Under without-replacement stratified sampling design and Hájek’s (1964) asymptotic framework, Berger and Torres (2016) showed that under  $H_0 : \boldsymbol{\theta}_U = \boldsymbol{\theta}_0$ , we have

$$\tilde{r}(\boldsymbol{\theta}, \boldsymbol{\varphi}_U) := 2 \left\{ \tilde{\ell}_{\max}(\boldsymbol{\varphi}_U) - \tilde{\ell}_{\max}(\boldsymbol{\theta}, \boldsymbol{\varphi}_U) \right\} \xrightarrow{d} \chi_{d_\theta}^2, \quad \text{if } \boldsymbol{\theta} = \boldsymbol{\theta}_0, \quad (9)$$

where

$$\begin{aligned} \tilde{\ell}_{\max}(\boldsymbol{\varphi}_U) &:= \max_{\mathbf{p}} \left\{ \tilde{\ell}(\mathbf{p}) : p_i > 0, \sum_{i \in \mathcal{S}} (p_i q_i - \psi_i) \frac{1}{\pi_i} \mathbf{f}(\mathbf{y}_i, \boldsymbol{\varphi}_U) = \mathbf{0}_{d_f}, \right. \\ &\quad \left. \sum_{i \in \mathcal{S}} (p_i q_i - \psi_i) \mathbf{z}_i = \frac{\mathbf{n}_H}{n} \right\}. \end{aligned}$$

In Berger and Torres (2016), the penalised empirical likelihood function uses the parametrization  $m_i := np_i \pi_i^{-1}$ . Straightforward algebra shows that it indeed reduces to (8). Note that the constraints within (8) do not imply  $\sum_{i \in \mathcal{S}} p_i = 1$ . We have  $\tilde{\ell}_{\max}(\boldsymbol{\theta}, \boldsymbol{\varphi}_U) = \ell_{\max}(\boldsymbol{\theta}, \boldsymbol{\varphi}_U)$  if we replace  $q_i$  by 1, because in this case,  $\psi_i = 0$  and  $\sum_{i \in \mathcal{S}} p_i = 1$ . In fact,  $q_i \rightarrow 1$ , as  $n/N \rightarrow 0$ . Thus, with  $n/N$  negligible,  $\tilde{\ell}_{\max}(\boldsymbol{\theta}, \boldsymbol{\varphi}_U) \approx \ell_{\max}(\boldsymbol{\theta}, \boldsymbol{\varphi}_U)$  and  $\tilde{r}(\boldsymbol{\theta}, \boldsymbol{\varphi}_U) \approx \hat{r}(\boldsymbol{\theta}^{(1)}, \boldsymbol{\varphi}_U)$ . Berger (2016) extended this §’s approach to Rao et al.’s (1962) sampling design with large sampling fraction.

## 4 Pseudoempirical likelihood

Wu and Rao's (2006) "pseudoempirical log-likelihood function" (see also Chen and Sitter, 1999) is defined by

$$L_{\max}(\boldsymbol{\theta}, \boldsymbol{\varphi}_{\mathcal{U}}) := \max_{\mathbf{p}} \left\{ L(\mathbf{p}) : p_i > 0, \sum_{i \in \mathcal{S}} p_i \mathbf{g}(\mathbf{y}_i, \boldsymbol{\theta}, \boldsymbol{\varphi}_{\mathcal{U}}) = \mathbf{0}_{d_g}, \sum_{i \in \mathcal{S}} p_i \mathbf{z}_i = \mathbf{1}_H \right\}, \quad (10)$$

where

$$\begin{aligned} L(\mathbf{p}) &:= n \sum_{i \in \mathcal{S}} \frac{\phi_i}{\pi_i} \log(p_i), \\ \phi_i &:= \frac{1}{N} \sum_{h=1}^H \frac{N_h}{\widehat{N}_h} z_{ih}, \\ N_h &:= \sum_{i \in \mathcal{U}} z_{ih}, \\ \widehat{N}_h &:= \sum_{i \in \mathcal{S}} \frac{z_{jh}}{\pi_j}. \end{aligned}$$

The function (10) clearly differs from (4). In (10),  $L(\mathbf{p})$  is different from Owen's (1988) empirical log-likelihood function  $\ell(\mathbf{p})$  and  $\sum_{i \in \mathcal{S}} p_i = H$ . In (4),  $\ell(\mathbf{p})$  is used and  $\sum_{i \in \mathcal{S}} p_i = 1$ . The main difference between (4) and (10) is the way in which the  $\pi_i$  are used. In (4), the  $\pi_i$  appears within the constraint. On the other hand, in (10), the  $\pi_i$  only appear within  $L(\mathbf{p})$ . The function  $L(\mathbf{p})$  is adjusted to take the  $\pi_i$  and the stratification into account by using  $\phi_i$ . The constraints involving the stratification variable  $\mathbf{z}_i$  also differ.

Wu and Rao (2006) showed that the "maximum pseudoempirical likelihood estimator", which maximises  $L_{\max}(\boldsymbol{\theta}, \boldsymbol{\varphi}_{\mathcal{U}})$ , is also the solution to

$$\widehat{\mathbf{E}}(\boldsymbol{\theta}, \boldsymbol{\varphi}_{\mathcal{U}}) := \sum_{i \in \mathcal{S}} \widehat{w}_i(\boldsymbol{\varphi}_{\mathcal{U}}) \mathbf{e}(\mathbf{y}_i, \boldsymbol{\theta}, \boldsymbol{\varphi}_{\mathcal{U}}) = \mathbf{0}_{d_g}, \quad (11)$$

where

$$\begin{aligned} \widehat{w}_i(\boldsymbol{\varphi}_{\mathcal{U}}) &:= n \widehat{\rho}_i(\boldsymbol{\varphi}_{\mathcal{U}}) \pi_i^{-1}, \\ \widehat{\rho}_i(\boldsymbol{\varphi}_{\mathcal{U}}) &:= n^{-1} \{1 + \boldsymbol{\lambda}(\boldsymbol{\varphi}_{\mathcal{U}})^\top \mathbf{u}_i(\boldsymbol{\varphi}_{\mathcal{U}})\}^{-1} \phi_i, \\ \mathbf{u}_i(\boldsymbol{\varphi}_{\mathcal{U}}) &:= \{\mathbf{f}(\mathbf{y}_i, \boldsymbol{\varphi}_{\mathcal{U}})^\top, \mathbf{z}_i^{(H-1)\top} - \mathbf{N}_{H-1} N^{-1}\}^\top, \\ \mathbf{z}_i^{(H-1)} &:= \{z_{i1}, \dots, z_{ih}, \dots, z_{i(H-1)}\}^\top, \\ \mathbf{N}_{H-1} &:= (N_1, \dots, N_h, \dots, N_{H-1})^\top. \end{aligned}$$

Here,  $\boldsymbol{\lambda}(\boldsymbol{\varphi}_{\mathcal{U}})$  is a Lagrangian parameter which is such that

$$\sum_{i \in \mathcal{S}} \widehat{w}_i(\boldsymbol{\varphi}_{\mathcal{U}}) \mathbf{u}_i(\boldsymbol{\varphi}_{\mathcal{U}}) = \mathbf{0}_{d_t+H-1}.$$

We see that the weights  $\widehat{w}_i(\boldsymbol{\varphi}_U)$  differ from  $\widehat{m}_i(\boldsymbol{\varphi}_U)$  given by (6), because  $\widehat{p}_i(\boldsymbol{\varphi}_U) \neq \widehat{p}_i(\boldsymbol{\varphi}_U)$  and  $\mathbf{u}_i(\boldsymbol{\varphi}_U) \neq \mathbf{c}_i(\boldsymbol{\varphi}_U)$ . Hence maximum pseudoempirical likelihood estimates are different from maximum empirical likelihood estimates. However, we expect minor numerical differences between them. The main difference between the approaches of §§ 4 and 3 is in the pivotal property of the empirical likelihood ratio statistics.

Suppose that  $d_{\boldsymbol{\theta}} = d_e = 1$ ; that is, we have a scalar parameter  $\theta_U$ . Let  $\widehat{\theta}$  be the maximum pseudoempirical likelihood estimator. Wu and Rao (2006) showed that under  $H_0 : \theta_U = \theta_0$ ,

$$\widehat{r}(\theta, \boldsymbol{\varphi}_U)_{\text{PEL}} := 2 \left\{ L_{\max}(\widehat{\theta}, \boldsymbol{\varphi}_U) - L_{\max}(\theta, \boldsymbol{\varphi}_U) \right\} \text{Deff}(\theta, \boldsymbol{\varphi}_U)^{-1} \xrightarrow{d} \chi_1^2, \quad \text{if } \theta = \theta_0,$$

where  $\text{Deff}(\theta, \boldsymbol{\varphi}_U)$  is called the “*design effect*” given by

$$\text{Deff}(\theta, \boldsymbol{\varphi}_U) := \frac{\text{Var}\{\widehat{E}(\theta, \boldsymbol{\varphi}_U)\}}{\text{Var}_{\text{SRS}}\{\widehat{E}(\theta, \boldsymbol{\varphi}_U)\}}. \quad (12)$$

Here,  $\widehat{E}(\theta, \boldsymbol{\varphi}_U)$  is defined by (11), when  $d_e = 1$ . The quantity  $\text{Var}\{\widehat{E}(\theta, \boldsymbol{\varphi}_U)\}$  denotes the variance under the sampling design and  $\text{Var}_{\text{SRS}}\{\widehat{E}(\theta, \boldsymbol{\varphi}_U)\}$  is the variance under simple random sampling ; that is, a sampling design with equal probability. The design effect would need to be estimated. We refer to Wu and Rao (2006) for more details about the estimation of (12).

The advantage of pseudoempirical likelihood is that it can be applied in principle to any complex sampling designs, because the design effect capture the complexity of any designs. The properties (7) and (9) are limited to single stage design or multi-stage design with small sampling fractions (Berger and Torres, 2016; Oğuz-Alper and Berger, 2016). The property (9) is valid under Hájek (1964) framework involving maximum entropy. Note that most designs used in practice have large entropy (Berger, 2011).

The disadvantage of pseudoempirical likelihood is that (12) only holds when the parameter is scalar ( $d_{\boldsymbol{\theta}} = d_e = 1$ ). For example, it cannot be used with multidimensional parameters, because the design effect implicitly assumes that (11) is scalar. Thus, for example, we do not a pivotal statistics for regression parameters. The properties (7) and (9) holds for multidimensional parameters. Thus, the approach of §3 can be used for generalised linear models. In (12), the design effect needs to be estimated, which adds some additional variability that may compromise to the convergence of  $\widehat{r}(\theta, \boldsymbol{\varphi}_U)_{\text{PEL}}$  towards the  $\chi^2$  distribution . Berger and Torres (2016) showed via a series of simulation that coverages of confidence intervals obtained from (7) and (9) are closer to the nominal value, than coverages obtained from (12). From a computational point of view (12) has the disadvantage of relying on variance estimates, which can be tedious to compute under complex sampling. Furthermore, (12) shows that the pseudoempirical likelihood ratio statistics is not pivotal, because of presence of a design effect.

## References

Berger, Y. G. (2011), “Asymptotic consistency under large entropy sampling designs with unequal probabilities,” *Pakistan Journal of Statistics*, 27(4), 407–426.

- Berger, Y. G. (2016), “Empirical Likelihood Inference for the Rao-Hartley-Cochran Sampling Design,” *Scandinavian Journal of Statistics*, 43, 721–735.
- Berger, Y. G., and Torres, O. D. L. R. (2016), “An empirical likelihood approach for inference under complex sampling design,” *Journal of the Royal Statistical Society Series B*, 78(2), 319–341.
- Chen, J., and Sitter, R. R. (1999), “A pseudo empirical likelihood approach to the effective use of auxiliary information in complex surveys,” *Statist. Sinica*, 9, 385–406.
- Chen, S., and Van Keilegom, I. (2009), “A review on empirical likelihood methods for regression,” *Test*, 18, 415–447.
- Deville, J. C., and Särndal, C.-E. (1992), “Calibration Estimators in Survey Sampling,” *Journal of the American Statistical Association*, 87(418), 376–382.
- Godambe, V. (1966), “A new approach to sampling from finite population I, II,” *Journal of the Royal Statistical Society Series B*, 28, 310–328.
- Godambe, V. P. (1960), “An Optimum Property of Regular Maximum Likelihood Estimation,” *The Annals of Mathematical Statistics*, 31(4), pp. 1208–1211.
- Hájek, J. (1964), “Asymptotic Theory of Rejective Sampling with Varying Probabilities from a Finite Population,” *The Annals of Mathematical Statistics*, 35(4), 1491–1523.
- Hartley, H. O., and Rao, J. N. K. (1968), “A new estimation theory for sample surveys,” *Biometrika*, 55(3), 547–557.
- Horvitz, D. G., and Thompson, D. J. (1952), “A Generalization of Sampling Without Replacement From a Finite Universe,” *Journal of the American Statistical Association*, 47(260), 663–685.
- Imbens, G. W., and Lancaster, T. (1994), “Combining Micro and Macro Data in Microeconomic Models,” *The Review of Economic Studies*, 61(4), 655–680.
- Kott, P. S. (2009), “Calibration weighting: combining probability samples and linear prediction models,” in *Sample Surveys: Design, Methods and Applications*, eds. D. Pfeffermann, and C. Rao, Handbook of Statistics, Amsterdam: Elsevier, pp. 55–82.
- Neyman, J. (1938), “On the Two Different Aspects of the Representative Method: The Method of Stratified Sampling and the Method of Purposive Selection,” *Journal of the Royal Statistical Society*, 97(4), 558–625.
- Oğuz-Alper, M., and Berger, Y. G. (2016), “Empirical likelihood approach for modelling survey data,” *Biometrika*, 103(2), 447–459.
- Owen, A. B. (1988), “Empirical Likelihood Ratio Confidence Intervals for a Single Functional,” *Biometrika*, 75(2), 237–249.
- Owen, A. B. (1991), “Empirical Likelihood for Linear Models,” *Ann. Statist.*, 19(4), 1725–1747.
- Owen, A. B. (2001), *Empirical Likelihood*, New York: Chapman & Hall.
- Polyak, B. T. (1987), *Introduction to Optimization*, New York: Optimization Software, Inc., Publications Division.
- Rao, J. N. K., Hartley, H. O., and Cochran, W. G. (1962), “On a Simple Procedure of Unequal Probability Sampling without Replacement,” *J. R. Stat. Soc. Ser. B. Stat. Methodol.*, 24(2), pp. 482–491.
- Wu, C., and Rao, J. N. K. (2006), “Pseudo-empirical likelihood ratio confidence intervals for complex surveys,” *Canad. J. Statist.*, 34(3), 359–375.