

# Design-based confidence intervals and significance test for regression parameters using an empirical likelihood approach

Melike Oguz-Alper\* ([M.OguzAlper@soton.ac.uk](mailto:M.OguzAlper@soton.ac.uk))<sup>1</sup>, Yves G. Berger ([Y.G.Berger@soton.ac.uk](mailto:Y.G.Berger@soton.ac.uk))<sup>1</sup>

**Keywords:** Design-based inference, estimating equations, nuisance parameter, unequal inclusion probabilities.

## 1. INTRODUCTION

Confidence intervals based on least squares may have poor coverages for regression parameters when the effect of sampling design is ignored. In addition, confidence intervals obtained from the standard design-based approaches [e.g. 1, 2, 3, 4] may not have the right coverages when the sampling distribution is skewed.

We propose to use an empirical likelihood approach to construct design-based confidence intervals and to test hypotheses for regression parameters under unequal probability sampling. Berger and De La Riva Torres [5] proposed an empirical likelihood approach which can be used for point estimation and to construct confidence intervals under complex sampling designs for a single parameter. We show that this approach can be extended to the multidimensional parameter case, in the sense that we can derive confidence intervals and test the significance of a subset of model parameters while taking the sampling design into account. This requires profiling which is not covered by Berger and De La Riva Torres [5].

The proposed approach intrinsically incorporates sampling weights, design variables, and auxiliary information. It may yield to more accurate confidence intervals when the sampling distribution of the regression parameters is not normal, the point estimator is biased, or the regression model is not linear. The proposed approach is simple to implement and less computer intensive than bootstrap. It does not rely on re-sampling, linearisation, variance estimation, or design-effect.

### 1.1. Parameter of interest and estimating equations

Let  $s$  be a random sample of size  $n$  which is selected from the finite population  $U$  of size  $N$  with respect to a probability sampling  $p(s)$ . Let  $y_i$  and  $\mathbf{x}_i$  be some variables of interest. Suppose that  $\psi_N$  is an unknown finite population parameter, which is the solution of the following population estimating equation.

$$\mathbf{G}(\psi) = \sum_{i \in U} \mathbf{g}_i(y_i, \mathbf{x}_i, \psi) = \mathbf{0},$$

where  $\mathbf{g}_i(y_i, \mathbf{x}_i, \psi)$  is a vector of estimating functions [e.g. 1, 2, 4, 6]. For example, for a simple linear regression, we have  $\mathbf{g}_i(y_i, \mathbf{x}_i, \psi) = \mathbf{x}_i(y_i - \mathbf{x}_i^\top \beta)$ .

We assume that the finite population parameter  $\psi_N$  converges to the model parameter  $\psi_0$ . If  $\hat{\psi}$  is a design-consistent estimator of  $\psi_N$  based on a sample data (see Section 2.1), the estimator  $\hat{\psi}$  is also an estimator of  $\psi_0$ . Assuming that the sampling fraction is negligible, the variability of  $\hat{\psi}$  is driven by the sampling design. Hence, design-based

---

<sup>1</sup> University of Southampton, Southampton Statistical Sciences Research Institute, Southampton, SO17 1BJ, United Kingdom. \* Funded by the Economic and Social Research Council (ESRC), United Kingdom.

confidence intervals proposed in this paper can be viewed as confidence intervals of  $\psi_N$  or  $\psi_0$ .

## 2. EMPIRICAL LIKELIHOOD INFERENCE

We use the *empirical log-likelihood function* given by Berger and De La Riva Torres [5]. It is defined as follows.

$$\ell(m) = \sum_{i \in s} \log(m_i), \quad (1)$$

where the  $m_i$  are unknown scale loads. The empirical log-likelihood function in (1) can be used for the sampling with replacement with unequal probability designs as shown by Hartley and Rao [7]. In this paper, we assume that the sampling fraction is negligible. Hence, the proposed approach is valid under the  $\pi_{\text{ps}}$  sampling as  $n/N \rightarrow 0$ .

The *maximum empirical likelihood estimators*  $\hat{m}_i$  maximise the empirical log-likelihood in (1) with respect to the constraints  $m_i \geq 0$  and

$$\sum_{i \in s} m_i \mathbf{c}_i = \mathbf{C}, \quad (2)$$

where the  $\mathbf{c}_i$  and  $\mathbf{C}$  are vectors defined in Section 2.1. We assume that  $\mathbf{c}_i$  and  $\mathbf{C}$  satisfy with a set of regularity conditions given by Berger and De La Riva Torres [5] and the condition  $\|\partial \mathbf{c}_i / \partial \lambda\| = O(1)$ , for all  $i \in s$  and  $\lambda \in \Lambda$ , where  $\|\cdot\|$  denotes the Euclidean norm,  $O(\cdot)$  defines the order of convergence, and  $\Lambda$  is a neighbourhood around the true population value  $\lambda_N$ . This condition implicitly implies that the  $\mathbf{c}_i$  are differentiable with respect to  $\lambda$  in a neighbourhood of  $\lambda_N$  [e.g. 1, 6, 8].

Berger and De La Riva Torres [5] showed that the maximum empirical likelihood estimators  $\hat{m}_i$  are given by  $\hat{m}_i = (\pi_i + \boldsymbol{\eta}^\top \mathbf{c}_i)^{-1}$ , where  $\boldsymbol{\eta}$  is such that the constraint (2) is satisfied.

### 2.1. Point estimation

Let  $\ell(\hat{m})$  be the maximum value of the empirical log-likelihood function  $\ell(m)$  under the constraints  $m_i \geq 0$  and (2) with  $c_i = \pi_i$  and  $C = n$ . This implies that  $\hat{m}_i = \pi_i^{-1}$ . Assume that  $\hat{m}_i^*$  maximises  $\ell(m)$  subject to the constraints  $m_i \geq 0$  and  $\sum_{i \in s} m_i \mathbf{c}_i^* = \mathbf{C}^*$  with  $\mathbf{c}_i^* = (c_i, \mathbf{g}_i(y_i, \mathbf{x}_i, \psi)^\top)^\top$  and  $\mathbf{C}^* = (C, \mathbf{0}^\top)^\top$ , for a given vector  $\psi$ . The *empirical log-likelihood ratio function* is defined by

$$\hat{r}(\psi) = 2\{\ell(\hat{m}) - \ell(\hat{m}^*(\psi))\}. \quad (3)$$

The *maximum empirical likelihood estimate*  $\hat{\psi}$  of the population parameter  $\psi_N$  is defined by the vector which minimises (3). The minimum value of (3) is obtained when  $\hat{r}(\psi) = 0$ ; that is, when  $\hat{m}_i^* = \hat{m}_i = \pi_i^{-1}$ . Thus, the maximum empirical likelihood estimator of  $\psi_N$  is the solution of the following sample estimating equation.

$$\hat{\mathbf{G}}(\psi) = \sum_{i \in s} \mathbf{g}_i(y_i, \mathbf{x}_i, \psi) \pi_i^{-1} = \mathbf{0}.$$

### 2.2. Hypothesis testing

Let  $\psi_N = (\boldsymbol{\theta}_N^\top, \boldsymbol{\lambda}_N^\top)^\top$  where  $\boldsymbol{\theta}_N$  is a  $p \times 1$  vector of parameters of interest and  $\boldsymbol{\lambda}_N$  is a

$q \times 1$  vector of parameters which are not of primary interest. Suppose we wish to test  $H_0 : \boldsymbol{\theta}_N = \boldsymbol{\theta}_N^0$ . Consider the *profile empirical log-likelihood ratio function* defined by

$$\widehat{r}(\boldsymbol{\theta}_N^0) = 2\{\ell(\widehat{m}) - \max_{\boldsymbol{\lambda}} \ell(\widehat{m}^*(\boldsymbol{\theta}_N^0, \boldsymbol{\lambda}))\}, \quad (4)$$

where the set of  $\widehat{m}_i^*$  maximises  $\ell(m)$  subject to the constraints  $m_i \geq 0$  and  $\sum_{i \in s} m_i \mathbf{c}_i^* = \mathbf{C}^*$  with  $\mathbf{c}_i^* = (\pi_i, \mathbf{g}_i(y_i, \mathbf{x}_i, \boldsymbol{\theta}_N^0, \boldsymbol{\lambda})^\top)^\top$  and  $\mathbf{C}^* = (n, \mathbf{0}^\top)^\top$ . Note that in (4), we maximise  $\ell(\widehat{m}^*(\boldsymbol{\theta}_N^0, \boldsymbol{\lambda}))$  over the parameter  $\boldsymbol{\lambda}$  for a given value of  $\boldsymbol{\theta}_N = \boldsymbol{\theta}_N^0$ .

Under  $H_0$ , it can be shown that the profile empirical log-likelihood ratio function  $\widehat{r}(\boldsymbol{\theta}_N^0)$  given by (4) follows asymptotically a *chi-squared distribution* with a  $p$  degree of freedom. Based on this, we can compute the *p-value*. Note that lack of fit would not affect the performance of the proposed empirical likelihood test [e.g. 8].

### 2.3. Confidence region

We can obtain confidence region for each parameter individually profiling out over the other parameters. In this case,  $p = 1$  and we have the scalar  $\theta_N$  instead of the vector  $\boldsymbol{\theta}_N$ . Then, based on the asymptotic chi-squared distribution of  $\widehat{r}(\theta_N^0)$  under the null hypothesis  $H_0 : \theta_N = \theta_N^0$ , the  $(1 - \alpha)\%$  empirical likelihood confidence region for  $\theta_N$  is given by the set  $\{\theta : \widehat{r}(\theta) \leq \chi_{df=1}^2(\alpha)\}$ , where  $\chi_{df=1}^2(\alpha)$  is the upper  $\alpha$  - *quantile* of the chi-squared distribution with one degree of freedom.

## 3. RESULTS

We present some numerical results for a linear regression model with one intercept and one slope. We generated the Hansen, Madow and Tepping (HMT) population [see 9]. The population size is  $N = 10\,000$ . We selected 1000 random samples of size  $n = 500$  from this population using the randomised systematic sampling with unequal probabilities.

The linear regression model of interest is defined by  $y_i = \lambda + \theta x_i + u_i$ , where the  $u_i$  are independent random variables with  $\text{var}(u_i|x_i) \propto x_i^{3/2}$ . The parameter of interest is the slope  $\theta$ . We profile out over the intercept  $\lambda$  when minimising (4).

Table 1 gives the observed coverages of the 95% confidence intervals constructed based on several methods. We considered two Pseudo likelihood approaches which are given by Binder and Patak [2] and Godambe and Thompson [4] [see also 1].

Standard confidence intervals are based on the normality of the point estimator. Note that, when the sampling distribution is skewed, the normality assumption may not hold. This explains the poor coverages of the Wald and the pseudo likelihood 1 approaches (see Table 1). The poor coverage of the Wald type of confidence intervals is also due to the fact that this method ignores the sampling design. We have an overcoverage with the rescaled bootstrap [e.g. 10]. Moreover, it has the largest confidence intervals on average compared to the other methods (see the ratio of average lengths in Table 1).

The coverage probabilities of the empirical likelihood and the pseudo likelihood 2 confidence intervals are not significantly different from the nominal level (i.e. 95%). However, the former is more reliable than the latter with regards to the standard deviation of length (see the last column of Table 1).

#### 4. CONCLUSIONS

We proposed an empirical likelihood approach which can be used to make inferences for regression parameters incorporating the sampling design. The proposed approach can be used for generalised linear models.

It can be easily shown that the population level information can be taken into account with the proposed approach. Unlike the usual calibration approach [11], the proposed approach can be used for testing and constructing confidence intervals. Moreover, the auxiliary information does not have to be in the form of totals or means [5].

The proposed approach can be easily extended to stratified sampling designs by incorporating the strata information into the  $c_i$ .

**Table 1. Observed coverages of the 95% confidence intervals for the slope  $\theta_N$ .**

N=10 000, n=500	Coverage probability	Lower error	Upper error	Ratio average length	Ratio SD length
Wald	76.6*	23.8*	0.1*	0.96	0.53
Empirical likelihood	94.8	3.1*	2.1*	1.00	1.00
Pseudo likelihood 1	94.0*	3.5*	2.5	0.97	1.07
Pseudo likelihood 2	94.8	3.3*	1.9*	0.99	1.09
Rescaled bootstrap	96.5*	2.4	1.1*	1.05	0.91

\* Significantly different from the nominal levels (95% and 2.5% for coverage probability and tail errors respectively) at the 5% significance level (i.e.  $p - value \leq 0.05$ ).

#### REFERENCES

- [1] D. A. Binder, On the variances of asymptotically normal estimators from complex surveys, *International Statistical Review*, 51, (1983), 279–292.
- [2] D. A. Binder and Z. Patak, Use of estimating functions for estimation from complex surveys, *Journal of the American Statisticsl Association*, 89, (1994), 1035–1043.
- [3] J. C. Deville, Variance estimation for complex statistics and estimators: linearization and residual techniques, *Survey Methodology*, 25, (1999), 193–203.
- [4] V. P. Godambe and M. E. Thompson, Estimating functions and survey sampling, *Handbook of Statistics: Design, Method and Applications*: D. Pfeffermann and C.R. Rao.(editors), Elsevier, 29B, (2009), 83–101.
- [5] Y. G. Berger and O. De La Riva Torres, *Empirical likelihood confidence intervals for complex sampling designs*. S3RI, <http://eprints.soton.ac.uk/337688>, (2012).
- [6] J. Qin and J. Lawless, Empirical likelihood and general estimating equations, *The Annals of Statistics*, 22, (1994), 300–325.
- [7] H. O. Hartley and J. N. K. Rao, *A new estimation theory for sample surveys*, II. Wiley-Interscience, New York, 1969.
- [8] A. B. Owen, *Empirical Likelihood*, Chapman & Hall, New York, 2001.
- [9] M. H. Hansen, W. G. Madow, and B. J. Tepping, An evaluation of model-dependent and probability-sampling inferences in sample surveys, *Journal of the American Statistical Association*, 78, (1983), 776–793.

- [10] J. N. K. Rao, C. F. J. Wu, and K. Yue, Some recent work on resampling methods for complex surveys, *Survey Methodology*, 18, (1992), 209–217.
- [11] J. C. Deville and C. E. Särndal, Calibration estimators in survey sampling, *Journal of the American Statistical Association*, 87, (1992), 376–382.