

# An R Library to construct empirical likelihood confidence intervals for complex estimators

Yves G. Berger, University of Southampton, UK

**Keywords:** Calibration, Design-based approach, Estimating equations, Finite population corrections, Hajek estimator, Horvitz-Thompson estimator, Regression estimator, Stratification, Unequal inclusion probabilities.

We developed an R library which can be used to compute empirical likelihood point estimates and confidence intervals. After explaining the empirical likelihood theory, we show how to use this library and an example based on the 2009 EU-SILC survey.

## 1. INTRODUCTION

Under complex sampling designs, point estimators may not have a normal sampling distribution and linearised variance estimators may be biased. Hence standard confidence intervals based upon the central limit theorem may have poor coverages. We propose an empirical likelihood approach which gives design based confidence intervals. The proposed approach does not rely on the normality of the point estimator, variance estimates, design-effects, re-sampling, joint- inclusion probabilities and linearisation, even when the estimator of interest is not linear. It can be used to construct confidence intervals for a large class of complex sampling designs and complex estimators which are solution of an estimating equation [4]. It can be used for means, regressions coefficients, quantiles, totals or counts even when the population size is unknown. It can be used with large and negligible sampling fractions. It also provides asymptotically optimal point estimators, and naturally includes calibration constraints [2]. The proposed approach is computationally simpler than the pseudo empirical likelihood [9] and the bootstrap approaches [8]. Berger and De La Riva Torres [1] show that the empirical likelihood confidence interval may give better coverages than the approaches based on linearisation [3], bootstrap [8] and pseudo empirical likelihood [9].

## 2. EMPIRICAL LIKELIHOOD APPROACH

Let  $U$  be a finite population of  $N$  units; where  $N$  is a fixed quantity which is not necessarily known. Suppose that the population parameter of interest  $\theta_0$  is the unique solution of the following estimating equation [4].

$$G(\theta) = 0, \quad \text{with } G(\theta) = \sum_{i \in U} g_i(\theta);$$

where  $g_i(\theta)$  is a function of  $\theta$  and of the characteristics of the unit  $i$ , such as the variables of interests and the auxiliary variables.

We propose to use the following *empirical log-likelihood function* [e.g. 1, 7].

$$\ell(m) = \sum_{i=1}^n \log(m_i),$$

where  $\sum_{i=1}^n$  denotes the sum over the sampled units. The quantities  $m_i$  are unknown positive scale loads. The maximum likelihood estimators of  $m_i$  are the values  $\hat{m}_i$  which maximise  $\ell(m)$  subject to the constraints  $m_i \geq 0$  and

$$\sum_{i=1}^n m_i \mathbf{c}_i = \mathbf{C};$$

where  $\mathbf{c}_i$  is  $Q \times 1$  vector associated with the  $i$ -th sampled unit and  $\mathbf{C}$  is  $Q \times 1$  vector (see Section 2.1). The values  $\hat{m}_i$  are survey weights.

## 2.1. Maximum empirical likelihood estimator

Suppose that the finite population  $U$  is stratified into  $H$  strata denoted by  $U_1, \dots, U_h, \dots, U_H$ ; where  $\cup_{h=1}^H U_h = U$ . Suppose that a sample  $s_h$  of fixed size  $n_h$  is selected with replacement with unequal probabilities from  $U_h$ . Let  $\mathbf{c}_i = \mathbf{z}_i$  and  $\mathbf{C} = \mathbf{n}$ ; where  $\mathbf{z}_i$  are the values of the design (or stratification) variables defined by  $\mathbf{z}_i = (z_{i1}, \dots, z_{iH})^\top$ , where  $\mathbf{n} = (n_1, \dots, n_H)^\top$  denotes the vector of the strata sample sizes, with  $z_{ih} = \pi_i$  when  $i \in U_h$  and  $z_{ih} = 0$  otherwise. It can be shown that  $\hat{m}_i = \pi_i^{-1}$ . Let  $\ell(\hat{m}) = \sum_{i=1}^n \log(\hat{m}_i)$  be the maximum value of the empirical log-likelihood function.

Let  $\hat{m}_i^*(\theta)$  be the values which maximise  $\ell(m)$  subject to the constraints  $m_i \geq 0$  and  $\sum_{i=1}^n m_i \mathbf{c}_i^* = \mathbf{C}^*$  with  $\mathbf{c}_i^* = (\mathbf{c}_i^\top, g_i(\theta))^\top$  and  $\mathbf{C}^* = (\mathbf{C}^\top, 0)^\top$ , for a given  $\theta$ . Let  $\ell(\hat{m}^*, \theta) = \sum_{i=1}^n \log(\hat{m}_i^*(\theta))$ . The *empirical log-likelihood ratio function* (or deviance) is defined by the following function of  $\theta$ .

$$\hat{r}(\theta) = 2 \{ \ell(\hat{m}) - \ell(\hat{m}^*, \theta) \}.$$

The *maximum empirical likelihood estimate*  $\hat{\theta}$  of  $\theta_0$  is defined by the value of  $\theta$  which minimises the function  $\hat{r}(\theta)$ . As the minimum value of  $\hat{r}(\theta)$  is zero,  $\hat{\theta}$  is the solution of  $\hat{r}(\theta) = 0$ . It can be easily shown that this implies that  $\hat{\theta}$  is the solution of the following estimating equation.

$$\hat{G}(\theta) = 0, \quad \text{with} \quad \hat{G}(\theta) = \sum_{i=1}^n \hat{m}_i g_i(\theta);$$

when  $g_i(\theta) = y_i - n^{-1}\theta\pi_i$ , we have  $\hat{G}(\theta) = \sum_{i=1}^n g_i(\theta)\pi_i^{-1}$  and  $\hat{\theta}$  is the Horvitz-Thompson estimator [6] given by  $\hat{Y}_\pi = \sum_{i=1}^n y_i \pi_i^{-1}$ . When  $g_i = y_i - \theta N^{-1}$ ,  $\hat{\theta}$  is the Hajek [5] ratio estimator  $\hat{Y}_H = N \hat{N}_\pi^{-1} \hat{Y}_\pi$ , where  $\hat{N}_\pi = \sum_{i=1}^n \pi_i^{-1}$ .

## 2.2. Empirical likelihood confidence intervals

Berger and De La Riva Torres [1] show that the random variable  $\hat{r}(\theta_0)$  follows asymptotically a chi-squared distribution with one degree of freedom. Thus, the  $\alpha$  level empirical likelihood confidence interval for the population parameter  $\theta_0$  is given by

$$\{ \theta : \hat{r}(\theta) \leq \chi_1^2(\alpha) \}.$$

Note that  $\hat{r}(\theta)$  is a convex non-symmetric function with a minimum when  $\theta$  is the maximum empirical likelihood estimator. This interval can be found using a bisection search method. This involves calculating  $\hat{r}(\theta)$  for several values of  $\theta$ . Berger and De La Riva Torres [1] showed how this approach can be used to accommodate large sampling fractions and non-response.

It is also possible to calibrate towards parameters more complex than totals. For example, we may want to calibrate with respect to population means, quantiles or variances. In this case, the calibration constraint is specified by the estimating equations  $\sum_{i=1}^n m_i \mathbf{f}_i(\mathbf{x}_i, \boldsymbol{\vartheta}_0) = \mathbf{0}$ ; where  $\mathbf{f}_i(\mathbf{x}_i, \boldsymbol{\vartheta}_0)$  is a vector function of the auxiliary variables and of a known parameter  $\boldsymbol{\vartheta}_0$  which is the solution of the following estimating equation

$$\sum_{i \in U} \mathbf{f}_i(\mathbf{x}_i, \boldsymbol{\vartheta}_0) = \mathbf{0}.$$

Calibration constraints are taken into account by including the auxiliary variables within the  $\mathbf{c}_i$ . In this case, we use  $\mathbf{c}_i = (\mathbf{z}_i^\top, \mathbf{f}_i(\mathbf{x}_i, \boldsymbol{\vartheta}_0)^\top)^\top$  and  $\mathbf{C} = (\mathbf{n}^\top, \mathbf{0}^\top)^\top$ . For example, if we want to calibrate towards known population means  $\boldsymbol{\vartheta}_0 = \mathbf{X}N^{-1}$ , we need to use  $\mathbf{f}_i(\mathbf{x}_i, \boldsymbol{\vartheta}_0) = \mathbf{x}_i - \boldsymbol{\vartheta}_0$ ; where  $\mathbf{X} = \sum_{i \in U} \mathbf{x}_i$  is a vector of known population totals. Simultaneous calibration on totals, means, proportion or any known parameter is also feasible.

### 3. AN R LIBRARY

In order to implement the approach describes in Section 2, we developed a library in R called *emplikpop*. First, we need to specify the design and calibration variables. Secondly, we need to specify the parameter of interest (i.e. the definition of  $g_i(\theta)$ ) and the finite population correction. This information is needed for point estimation and for confidence intervals.

- i. *Design and auxiliary information:* This information is contained in the matrix `MatDA` and the vector `Var.Labels`. `MatDA` is a  $n \times p$  matrix containing the stratification labels, the  $\pi_i$  and  $\mathbf{f}_i(\mathbf{x}_i, \boldsymbol{\vartheta}_0)$  (depending on  $\boldsymbol{\vartheta}_0$ ). The vector `Var.Labels` gives the columns (in `MatDA`) of stratification and  $\pi_i$ .
- ii. *Definition of the parameter of interest:* A function object `FunctG` defining the function  $g_i(\theta)$ . This function depends on a matrix `Data` and a vector `Vect.Const` which specify the data needed in the definition of  $g_i(\theta)$ . This function also depends on `Theta`, a given value for  $\theta$ . This function has the following format:  
`FunctG = function(Data, Vect.Const, Theta){...}`
- iii. *Finite population corrections:* The  $n \times 1$  vector `fpc` contains the  $1 - \pi_i$  or  $1 - nN^{-1}$  or 1 (if the finite population correction is ignored).
- iv. *Survey weights:* The survey weight  $\hat{m}_i$  are computed using the function  
`> Mi = Vect.Mi(MatDA, Var.Labels)`
- v. *Point estimate:* The point estimate is computed using the following function  
`> Theta.Hat = SolEE(FunctG, Data, Vect.Const, Mi, Min, Max)`  
 where `Min` and `Max` specify the range of  $\hat{\theta}$ .
- vi. *Confidence interval:* The confidence interval is computed using the following function  
`> Bounds = ELBound(MatDA, Var.Labels, FunctG,  
 Level, Data, Vect.Const, Theta.Hat, fpc)`  
 where `Level` is the level of the confidence interval. For example, for the 95% confidence interval, we use `Level = 0.95`. The object `Bounds` contains the bounds of the confidence interval. We have also predefined function for means, totals and quantiles: `ELBoundMean()`, `ELBoundTotal()` and `ELBoundQuantile()`. For the Rao-Hartley-Cochran sampling design, we have the function `ELBound.RHC()`.

#### 4. AN APPLICATION TO THE EU-SILC HOUSEHOLD SURVEY

We use the 2009 EU-SILC user database to estimate the *persistent at-risk-of-poverty rate*. we adopted an ultimate cluster approach, where the units are the primary sampling units. In the table below, we have the point estimate and several confidence intervals for a couple of countries: the empirical likelihood confidence intervals, the standard confidence intervals based on variance estimates and the rescaled bootstrap confidences intervals [8]. Note that the bounds of the standard intervals are negative for Ireland, Austria, Malta and Denmark. The bootstrap bounds and the empirical likelihood bounds are larger than the bounds of the standard intervals. These differences are more pronounced for Austria, Malta, Denmark, the Netherlands, Estonia, Latvia and Greece. This is due to the skewness of the sampling distribution.

**Table 1: Persistent at-risk-of-poverty rate & confidence intervals. 2009 EU-SILC.**

Country	Rate (%)	Emp. Likelihood		Standard		Rescaled Bootstrap	
		Lower	Upper	Lower	Upper	Lower	Upper
Ireland	0.53	0.08	1.76	-0.26	1.31	0.00	1.58
Austria	2.14	0.53	6.50	-0.52	4.80	0.14	5.26
Malta	2.90	0.97	7.75	-0.10	5.89	0.62	6.09
Denmark	3.46	1.09	8.95	-0.06	6.98	0.67	7.76
France	4.50	3.33	5.99	3.21	5.8	3.23	6.04
UK	5.18	2.56	9.90	1.78	8.57	2.15	8.85
Netherlands	5.22	1.88	11.66	0.69	9.75	1.31	10.25
Estonia	7.45	4.07	14.69	2.87	12.03	3.47	13.11
Poland	8.58	5.89	12.49	6.32	10.85	5.32	12.13
Latvia	10.34	6.09	17.36	5.05	15.63	5.36	15.27
Greece	11.34	7.51	18.32	6.72	15.96	7.03	16.95

#### REFERENCES

- [1] Berger, Y. G. and De La Riva Torres, O. (2012) Empirical likelihood confidence intervals for complex sampling designs. S3RI Working paper, <http://eprints.soton.ac.uk/337688/>
- [2] Deville, J. C. and Sarndal (1992) Calibration estimators in survey sampling. Journal of the American Statistical Association, 87, 376-382.
- [3] Deville, J. C. (1999) Variance estimation for complex statistics and estimators: linearization and residual techniques. Survey Methodology, 25} 193-203.
- [4] Godambe, V. P. (1960) An optimum property of regular maximum likelihood estimation. The Annals of Mathematical Statistics, 31, 1208-1211.
- [5] Hajek, J. (1971) Comment on a paper by D. Basu. in Foundations of Statistical Inference. Toronto: Holt, Rinehart and Winston.
- [6] 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, 663-685.
- [7] Owen, A. B. (2001) Empirical Likelihood. New York: Chapman & Hall.
- [8] Rao, J. N. K., Wu, C. F. J. and Yue, K. (1992) Some recent work on resampling methods for complex surveys. Survey Methodology, 18, 209--217.

[9] Wu, C. and Rao, J. N. K. (2006) Pseudo-empirical likelihood ratio confidence intervals for complex surveys. *The Canadian Journal of Statistics*, 34, 359-375.