

Empirical Likelihood confidence regions for regression parameters under unequal probability sampling

Yves G.BERGER*

University of Southampton

Southampton Statistical Sciences Research Institute

Southampton, SO17 1BJ, United Kingdom

Y.G.Berger@soton.ac.uk

We propose a new empirical likelihood approach which can be used to construct design-based confidence regions of regression parameters under unequal probability sampling. The proposed approach gives confidence regions which may have better coverages than standard confidence regions and pseudo empirical likelihood confidence regions which rely on variance estimates and design-effects. The proposed approach does not rely on variance estimates, design-effects, re-sampling or linearisation, even when the regression parameter is not linear. It also gives suitable confidence regions when the point estimator is biased. The proposed approach can also be adjusted to account for large sampling fractions.

Key Words: Design-based approach, Estimating equations, Regression estimator, Unequal inclusion probabilities.

Introduction

Let U be a finite population of N units; where N is a fixed quantity which is not necessarily known. Let θ_0 be an r -dimensional population parameter of interest. Suppose that θ_0 is the unique solution of the following set of estimating equation (e.g. Qin and Lawless, 1994).

$$(1) \quad \mathbf{G}(\theta) = \mathbf{0}, \quad \text{with} \quad \mathbf{G}(\theta) = \sum_{i \in U} \mathbf{g}_i(\theta);$$

where $\mathbf{g}_i(\theta)$ is an r -dimensional function of θ and of characteristics of the unit i . This function does not need to be differentiable.

Suppose we are interested in modelling the relationship between a response (or dependent) variable y_i and some r -dimensional covariates (or independent variables) \mathbf{u}_i . For example, with a linear model, θ_0 is a regression parameter and

$$\mathbf{g}_i(\theta) = \mathbf{u}_i(y_i - \mathbf{u}_i'\theta).$$

With a non-linear regression, we have

$$\mathbf{g}_i(\theta) = \frac{\partial f(\mathbf{u}_i, \theta)}{\partial \theta} \{y_i - f(\mathbf{u}_i, \theta)\},$$

where $f(\mathbf{u}_i, \theta)$ is a non linear function. For a robust regression, we have

$$\mathbf{g}_i(\theta) = \mathbf{u}_i\psi(y_i, \mathbf{u}_i'\theta),$$

where $\psi(\cdot)$ is a robust function (Huber, 1981).

Suppose that we wish to estimate θ_0 from the data of a sample s of size n selected with a single stage unequal probabilities without replacement sampling design. We consider that the sample size n is a fixed (non random) quantity. We adopt a design-based (non-parametric) approach; where the sampling distribution is specified by the sampling design and where the values of the response variable and the covariates are fixed (non random) quantities.

Empirical likelihood approach

We propose to use the following empirical likelihood function (e.g. Owen, 2001).

$$(2) \quad L(m) = \prod_{i=1}^n \frac{m_i}{N},$$

where m_i is the unit mass of unit i in the population (e.g. Deville, 1999).

Hartley and Rao (1969) showed that (2) is the an empirical likelihood function under unequal probability sampling with replacement, as m_i/N is the probability to observe the i -th unit. Owen (2001, Ch. 6) showed that (2) is a suitable empirical likelihood function when the units are selected independently with a Poisson sampling design. Although under fixed size sampling designs, the units are not selected independently, we propose to use the empirical likelihood function (2) under fixed size sampling designs. The aim is to show that this empirical likelihood function can be used for point estimation and to construct confidence regions (or to derive tests) under fixed size sampling designs.

The maximum likelihood estimators of m_i are the values \hat{m}_i which maximise the *log-empirical likelihood function*

$$(3) \quad \ell(m) = \sum_{i=1}^n \log(m_i),$$

subject to the constraints $m_i \geq 0$ and

$$(4) \quad \sum_{i=1}^n m_i \mathbf{c}_i = \mathbf{C};$$

where $\sum_{i=1}^n$ denotes the sum over the sampled units, \mathbf{c}_i is a known $Q \times 1$ vector associated with the i -th sampled unit and \mathbf{C} is a known $Q \times 1$ vector. We also assume that the constraint (4) is such that the fixed size constraint

$$(5) \quad \sum_{i=1}^n m_i \pi_i = n$$

always holds, where π_i denotes the inclusion probability of unit i . Under equal probability sampling, we have that $\pi_i = n/N$, and the constraint (5) reduces to $\sum_{i=1}^n m_i = N$ which is the constraint adopted under equal probability sampling (e.g. Rao and Wu, 2009). Berger and De La Riva Torres (2012) showed that the solution of this maximisation is given by

$$(6) \quad \hat{m}_i = (\pi_i + \boldsymbol{\eta}' \mathbf{c}_i)^{-1},$$

The quantity $\boldsymbol{\eta}$ is such that the constraint (4) holds. This quantity can be computed using an iterative Newton-Raphson procedure (e.g. Berger and De La Riva Torres, 2012; Rao and Wu, 2009).

The *maximum empirical likelihood estimator* $\hat{\boldsymbol{\theta}}$ of $\boldsymbol{\theta}_0$ is defined by solution of the following estimating equation.

$$(7) \quad \hat{\mathbf{G}}(\boldsymbol{\theta}) = \mathbf{0}, \quad \text{with} \quad \hat{\mathbf{G}}(\boldsymbol{\theta}) = \sum_{i=1}^n \hat{m}_i \mathbf{g}_i(\boldsymbol{\theta});$$

where \hat{m}_i is defined by (6). We assume that the $\mathbf{g}_i(\boldsymbol{\theta})$ are such that $\hat{\mathbf{G}}(\boldsymbol{\theta}) = \mathbf{0}$ has a unique solution. The estimator $\hat{\boldsymbol{\theta}}$ is a maximum empirical likelihood estimator because it also minimises the empirical log-likelihood ratio function (or deviance) defined by (8).

For example, suppose that we would like to fit a linear model with only an intercept ($\mathbf{u}_i = 1$). In this case, $\mathbf{g}_i(\boldsymbol{\theta}) = y_i - \theta$. If $\mathbf{c}_i = \pi_i$ and $\mathbf{C} = n$, we have that $m_i = \pi_i^{-1}$ and

$$\hat{\mathbf{G}}(\boldsymbol{\theta}) = \sum_{i=1}^n \frac{g_i(\boldsymbol{\theta})}{\pi_i}$$

is the Horvitz and Thompson (1952) estimator of the function (1). As $\mathbf{g}_i(\boldsymbol{\theta}) = y_i - \theta$, we have that $\hat{\theta}$ is the Hájek (1971) estimator of the population mean; that is, $\hat{\theta} = \hat{N}_\pi^{-1} \sum_{i=1}^n y_i \pi_i^{-1}$, where $\hat{N}_\pi = \sum_{i=1}^n \pi_i^{-1}$.

Empirical log-likelihood ratio function

The main advantage of the empirical likelihood approach is its capability of deriving non-parametric confidence regions which do not depend on variance estimates or on the normality of the point estimator. In this Section, we propose to use the empirical log-likelihood ratio function defined by (8) to derive empirical likelihood confidence regions.

Let \hat{m}_i be the values which maximise (3) subject to the constraints $m_i \geq 0$ and (4) when $\mathbf{c}_i = \pi_i$ and $\mathbf{C} = n$. Note that $m_i = \pi_i^{-1}$ in this situation. The maximum value of the empirical log-likelihood function is given by $\ell(\hat{m}) = -\sum_{i=1}^n \log(\pi_i)$. Let \hat{m}_i^* be the values which maximise (3) subject to the constraints $m_i \geq 0$ and (4) with $\mathbf{c}_i = \mathbf{c}_i^*$ and $\mathbf{C} = \mathbf{C}^*$, where $\mathbf{c}_i^* = (\pi_i, \mathbf{g}_i(\boldsymbol{\theta}))'$ and $\mathbf{C}^* = (n, \mathbf{0})'$. Let $\ell(\hat{m}^*, \boldsymbol{\theta})$ be the maximum value of the empirical log-likelihood function. The *empirical log-likelihood ratio function* (or deviance) is defined by the following function of $\boldsymbol{\theta}$.

$$(8) \quad \hat{r}(\boldsymbol{\theta}) = 2 \{ \ell(\hat{m}) - \ell(\hat{m}^*, \boldsymbol{\theta}) \}.$$

It can be easily shown that $\hat{r}(\hat{\boldsymbol{\theta}}) = \mathbf{0}$. Hence $\hat{\boldsymbol{\theta}}$ is indeed the maximum empirical likelihood estimator of $\boldsymbol{\theta}_0$, because it minimises the empirical log-likelihood ratio function. Berger and De La Riva Torres (2012) showed how the stratification can be taken into account by including the stratification variables within the vectors \mathbf{c}_i and \mathbf{c}_i^* .

We will show that under a set of regularity conditions, $\hat{r}(\boldsymbol{\theta}_0)$ follows asymptotically a chi-squared distribution with r degree of freedom when the sampling fraction, n/N , is negligible. This property relies on the fact that

$$(9) \quad \hat{\mathbf{G}}_\pi(\boldsymbol{\theta}_0) = \sum_{i=1}^n \frac{\mathbf{g}_i(\boldsymbol{\theta}_0)}{\pi_i},$$

is a vector of Horvitz and Thompson (1952) estimator which follows a multivariate normal distribution asymptotically (Berger, 1998; Hájek, 1964; Vísek, 1979).

Empirical likelihood confidence regions

As $\hat{r}(\boldsymbol{\theta}_0)$ follows asymptotically a chi-squared distribution, the $(1 - \alpha)$ level empirical likelihood confidence region (Wilks, 1938) for the population parameter $\boldsymbol{\theta}_0$ is given by the following set

$$(10) \quad \{ \boldsymbol{\theta} : \hat{r}(\boldsymbol{\theta}) \leq \chi_{df=p}^2(\alpha) \},$$

where $\chi_{df=p}^2(\alpha)$ is the upper α -quantile of the chi-squared distribution with r degree of freedom. Note that $\hat{r}(\boldsymbol{\theta})$ is a convex non-symmetric function with a minimum when $\boldsymbol{\theta}$ is the maximum empirical likelihood estimator. This region can be found using a bijection search method. This involves calculating $\hat{r}(\boldsymbol{\theta})$ for several values of $\boldsymbol{\theta}$.

Let $\boldsymbol{\theta}_0 = (\boldsymbol{\theta}'_1, \boldsymbol{\theta}'_2)'$; where $\boldsymbol{\theta}_1$ is an r_1 -dimensional vector ($r_1 < r$). Suppose that we would like to derive a confidence region for the parameter $\boldsymbol{\theta}_1$. The parameter $\boldsymbol{\theta}_2$ is treated as a nuisance parameter. Consider

$$\hat{r}(\boldsymbol{\theta}_1) = 2 \{ \ell(\hat{m}) - \text{Max}_{\boldsymbol{\theta}_2} \ell(\hat{m}^*, \boldsymbol{\theta}) \};$$

where $\text{Max}_{\boldsymbol{\theta}_2} \ell(\hat{m}^*, \boldsymbol{\theta})$ is the maximum value of $\ell(\hat{m}^*, \boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}_2$, for a fixed value of $\boldsymbol{\theta}_1$. We will show that $\hat{r}(\boldsymbol{\theta}_1)$ follows asymptotically a chi-squared distribution with r_1 degree of freedom when

the sampling fraction is negligible. Hence, the set $\{\theta_1^\circ : \hat{r}(\theta_1^\circ) \leq \chi_{df=r_1}^2(\alpha)\}$ gives a confidence region for θ_1 .

Empirical Likelihood Approach with Auxiliary Variables

In practice, population control totals of auxiliary variables are often known and this information is often taken into account at the estimation stage. Let x_i be the value of an auxiliary variables attached to unit i . Suppose that the control total $X = \sum_{i \in U} x_i$ is known. Traditional approaches (Birch, 1963) consist in adding these auxiliary variables within the covariates \mathbf{u}_i . However, in this situation, a different model is fitted and in practice we may want not to include the variable x_i within \mathbf{u}_i . In the approach proposed, we do not include x_i within \mathbf{u}_i . The variable x_i is use to derive the weights \hat{m}_i of the estimating equation.

Let $\mathbf{c}_i = (x_i, \pi_i)'$ and $\mathbf{C} = (X, n)'$. Berger and De La Riva Torres (2012) showed that

$$(11) \quad \hat{\mathbf{G}}(\boldsymbol{\theta}) = \hat{\mathbf{G}}_\pi(\boldsymbol{\theta}) + \hat{\mathbf{B}}(X - \hat{X}_\pi) + \hat{\boldsymbol{\epsilon}},$$

where $\hat{X}_\pi = \sum_{i=1}^n \check{x}_i$, $\|\hat{\boldsymbol{\epsilon}}\| = o_p(N)$,

$$\hat{\mathbf{B}} = \frac{\sum_{i=1}^n (\check{x}_i - n^{-1}\hat{X}_\pi)(\check{\mathbf{g}}_i(\boldsymbol{\theta}) - n^{-1}\hat{\mathbf{G}}_\pi(\boldsymbol{\theta}))}{\sum_{i=1}^n (\check{x}_i - n^{-1}\hat{X}_\pi)^2},$$

$\check{x}_i = x_i \pi_i^{-1}$ and $\check{\mathbf{g}}_i(\boldsymbol{\theta}) = \mathbf{g}_i(\boldsymbol{\theta}) \pi_i^{-1}$. Note that $\hat{\mathbf{B}}$ is the estimator of the covariance between $\hat{\mathbf{G}}_\pi(\boldsymbol{\theta})$ and \hat{X}_π divided by the estimator of the variance of \hat{X}_π under a with replacement pps sampling design (e.g. Särndal et al., 1992, p. 99). Therefore $\hat{\mathbf{B}}$ is the optimal regression coefficient (e.g. Berger et al., 2003; Isaki and Fuller, 1982; Montanari, 1987; Rao, 1994; Särndal, 1996) when the sampling fraction is negligible. Hence, the empirical likelihood estimator is asymptotically optimal. The result (11) can be generalised when \mathbf{x}_i is a vector of values of auxiliary variables. When N is known, we recommend to use $x_i = 1$ or to include a variable equal to one into \mathbf{x}_i . This may improve the efficiency of the maximum empirical likelihood estimator.

With auxiliary variables, the confidence regions have to be constructed using the following restricted empirical likelihood approach proposed by (Berger and De La Riva Torres, 2012), because the function (8) may not converge to a chi-square distribution. Let $\mathbf{c}_i = \dot{\mathbf{c}}_i$, $\mathbf{c}_i^* = (\dot{\mathbf{c}}_i', \mathbf{g}_i(\boldsymbol{\theta}_0))'$, $\mathbf{C} = (\mathbf{X}', \mathbf{n}')'$, and $\mathbf{C}^* = (\mathbf{X}', \mathbf{n}', \mathbf{0})'$, with $\dot{\mathbf{c}}_i = (\mathbf{x}'_i, \tilde{\pi}_i(\mathbf{x}))'$; where $\tilde{\pi}_i(\mathbf{x}) = \hat{m}_i(\mathbf{x})^{-1}$. Let $\ell(\hat{m}, \mathbf{x})$ be the maximum value (3) under the constraint (4) with \mathbf{c}_i and \mathbf{C} . Let $\ell(\hat{m}^*, \mathbf{x}, \boldsymbol{\theta})$ be the maximum value (3) under the constraint (4) with \mathbf{c}_i^* and \mathbf{C}^* . In both cases, we consider that the constraints (4) are such that $\sum_{i=1}^n m_i \tilde{\pi}_i(\mathbf{x}) = n$. The *restricted empirical log-likelihood ratio function* is given by

$$(12) \quad \hat{r}_x(\boldsymbol{\theta}) = 2 \{ \ell(\hat{m}, \mathbf{x}) - \ell(\hat{m}^*, \mathbf{x}, \boldsymbol{\theta}) \}.$$

Berger and De La Riva Torres (2012) showed that $\hat{r}_x(\boldsymbol{\theta}_0)$ follows asymptotically a chi-squared distribution with r degree of freedom. Note that $\tilde{\pi}_i(\mathbf{x})$ can be replaced by the inverse of any calibration weights (e.g. Deville and Särndal, 1992) which are calibrated with respect to $(\mathbf{x}'_i, \pi_i)'$, as the restricted empirical log-likelihood ratio function still follows a chi-squared distribution in this situation. Note that $\tilde{\pi}_i(\mathbf{x})$ may be larger than one. However, this does not cause any problem. Berger and De La Riva Torres (2012) showed how the stratification can be taken into account by including the stratification variables within the vectors \mathbf{c}_i and \mathbf{c}_i^* .

Non-negligible sampling fractions

With large sampling fractions, the empirical log-likelihood ratio function does not necessarily follow a chi-squared distribution. Berger and De La Riva Torres (2012) proposed to adjust the constraint in

order to obtain a chi-squared distribution asymptotically. Consider $\mathbf{c}_i = \pi_i$ and $\mathbf{C} = n$. We propose to use $\mathbf{c}_i^* = q_i(\pi_i, \mathbf{g}_i(\boldsymbol{\theta}))'$ and $\mathbf{C}^* = (\sum_{i=1}^n q_i, \sum_{i=1}^n (q_i - 1)\mathbf{g}_i(\boldsymbol{\theta})\pi_i^{-1})'$, with $q_i = (1 - \pi_i)^{1/2}$. Let \widehat{m}_i^* be defined by

$$(13) \quad \widehat{m}_i^* = \left(\pi_i + \boldsymbol{\eta}^{*'} \mathbf{c}_i^* \right)^{-1},$$

where $\boldsymbol{\eta}^*$ is such that $\sum_{i=1}^n \widehat{m}_i^* \mathbf{c}_i^* = \mathbf{C}^*$ holds. We propose to use the same empirical log-likelihood ratio function (8). The empirical log-likelihood ratio function is still defined by (8) with $\ell(m)$ given by (3). We will show that under a set of regularity conditions, $\widehat{r}(\boldsymbol{\theta}_0)$ follows asymptotically a chi-squared distribution with r degree of freedom for any sampling fractions. Hence empirical likelihood confidence regions can be constructed using (10). Berger and De La Riva Torres (2012) showed how the stratification can be taken into account by including the stratification variables within the vectors \mathbf{c}_i and \mathbf{c}_i^* .

The q_i are finite population corrections factors proposed by Berger (2005). The q_i reduce the effect on the confidence region of units with large π_i . For example, if $\pi_i = 1$, then $\widehat{m}_i \pi_i = \widehat{m}_i^* \pi_i = 1$. This implies that this unit will have no contribution towards the empirical likelihood functions and any confidence regions. This is a natural property as this unit does not contribute towards the sampling distribution. Note that we propose to adjust the constraints by quantities which do not need to be estimated, unlike the pseudo empirical likelihood approach (Wu and Rao, 2006) which adjusts the empirical log-likelihood ratio function by a quantity that needs to be estimated (the design effect).

Conclusions

Standard confidence regions based upon the central limit theorem can perform poorly when the sampling distribution is not normal. For example, the lower bounds of a confidence region can be negative even when the parameter of interest is positive. The coverage and the tail errors can be also different from their intended levels. On the other hand, empirical likelihood confidence regions may be better in this situation, as empirical likelihood confidence regions are determined by the distribution of the data (Rao and Wu, 2009) and as the range of the parameter space is preserved. Note that the distribution of a point estimator of is not necessarily normal, and the proposed empirical likelihood approach does not rely on the normality of the point estimator.

Standard confidence regions based on the central limit theorem require normality and variance estimates which often involve linearisation or re-sampling. The proposed method does not rely on normality, variance estimates, linearisation or re-sampling, even if the parameter of interest is not linear. Empirical likelihood confidence regions can be easier to compute than standard confidence regions based on variance estimates. It provides an alternative to bootstrap, when linearisation cannot be used. The proposed approach has some advantages over the bootstrap approach. It is less computationally intensive than the bootstrap. The proposed approach naturally includes auxiliary information and the information about the sampling design. This is particularly useful under informative sampling (Pfeffermann, 2009).

REFERENCES

Berger, Y. G. (1998), "Rate of Convergence to Normal Distribution for the Horvitz-Thompson Estimator," *Journal of Statistical Planning and Inference*, 67, 209–226.
 Berger, Y. G. (2005), "Variance Estimation with Highly Stratified Sampling Designs with Unequal Probabilities," *Australian and New Zealand Journal of Statistics*, 47(3), 365–373.
 Berger, Y. G., and De La Riva Torres, O. (2012), "A unified theory of empirical likelihood ratio confidence intervals for survey data with unequal probabilities and non negligible sampling fractions," *Southampton Statistical Sciences Research Institute*, <http://eprints.soton.ac.uk/337688>.

- Berger, Y. G., Tirari, M. E. H., and Tillé, Y. (2003), "Towards Optimal Regression Estimation in Sample Surveys," *Australian and New Zealand Journal of Statistics*, 45, 319–329.
- Birch, M. W. (1963), "Maximum likelihood in three-way contingency tables," *Journal of the Royal Statistical Society*, B(25), 220–233.
- Deville, J. C. (1999), "Variance estimation for complex statistics and estimators: linearization and residual techniques," *Survey Methodology*, 25, 193–203.
- Deville, J. C., and Särndal, C. E. (1992), "Calibration Estimators in Survey Sampling," *Journal of the American Statistical Association*, 87(418), 376–382.
- Hájek, J. (1964), "Asymptotic Theory of Rejective Sampling with Varying Probabilities from a Finite Population," *The Annals of Mathematical Statistics*, 35(4), 1491–1523.
- Hájek, J. (1971), "Comment on a paper by D. Basu. In Foundations of Statistical Inference. Toronto: Holt, Rinehart and Winston,".
- Hartley, H. O., and Rao, J. N. K. (1969), *A new estimation theory for sample surveys, II*, A Symposium on the Foundations of Survey Sampling held at the University of North Carolina, Chapel Hill, North Carolina: Wiley-Interscience, New York.
- 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.
- Huber, P. J. (1981), *Robust Statistics* Wiley, New York.
- Isaki, C. T., and Fuller, W. A. (1982), "Survey design under the regression super-population model," *Journal of the American Statistical Association*, 77, 89–96.
- Montanari, G. (1987), "Post sampling efficient QR-prediction in large sample survey," *International Statistical-Review*, 55, 191–202.
- Owen, A. B. (2001), *Empirical Likelihood*, New York: Chapman & Hall.
- Pfeffermann, D. (2009), "Inference under informative sampling," *Handbook of Statistics 29: Sample Surveys: Inference and Analysis*. Elsevier. Danny Pfeffermann and C.R. Rao eds, pp. 455–487.
- Qin, J., and Lawless, J. (1994), "Empirical Likelihood and General Estimating Equations," *The Annals of Statistics*, 22(1), pp. 300–325.
- Rao, J. N. K. (1994), "Estimating total and distribution function using auxiliary information at the estimation stage," *Journal of Official Statistics*, 10(2), 153–165.
- Rao, J. N. K., and Wu, W. (2009), "Empirical Likelihood Methods," *Handbook of statistics: Sample Surveys: Inference and Analysis*, D. Pfeffermann and C. R. Rao eds. *The Netherlands (North-Holland)*, 29B, 189–207.
- Särndal, C. E. (1996), "Efficient estimators with simple variance in unequal probability sampling," *Journal of the American Statistical Association*, 91, 1289–1300.
- Särndal, C.-E., Swensson, B., and Wretman, J. (1992), *Model Assisted Survey Sampling*, New York: Springer-Verlag.
- Víšek, J. (1979), "Asymptotic distribution of simple estimate for rejectif, sampford and successive sampling," *Contribution to Statistics, Jaroslav Hajek Memorial Volume. Academia of Prague, Czech Republic*, pp. 71–78.
- Wilks, S. S. (1938), "Shortest Average Confidence Intervals from Large Samples," *The Annals of Mathematical Statistics*, 9(3), 166–175.
- Wu, C., and Rao, J. N. K. (2006), "Pseudo-empirical likelihood ratio confidence intervals for complex surveys," *The Canadian Journal of Statistics*, 34(3), 359–375.