

Large sample distribution of the sample total in a generalized rejective sampling scheme

S. Rao Jammalamadaka, G. Michaletzky * and P. Todorovic

Statistics and Applied Probability Program, University of California, Santa Barbara, CA 93106, USA

Received March 1990

Abstract: The weak convergence of a sample sum, in a generalized rejective sampling from a finite population, to a Poisson and Normal distribution is discussed. The generalization consists in assuming that the elements of the population are random variables, rather than fixed values.

1. Introduction

Consider a finite population of N units, $\{Y_j\}_1^N$, where the value of the j th unit, Y_j , is a non-negative integer. From this, a sample of size n is drawn according to a rejective sampling plan with parameters p_1, \dots, p_N , with $\sum_1^N p_j = n$ (see Hajek, 1981, Chapter 7). Here p_j denotes the probability that the j th population unit is included in the sample — this event being represented by the indicator variable, I_j . Let S_{N_n} be the sample sum, then clearly $\mathcal{L}\{S_{N_n}\} = \mathcal{L}\{\sum_1^N Y_j I_j \mid \sum_1^N I_j = n\}$, where $\{I_j\}_1^N$ is a sequence of independent Bernoulli r.v.'s with $E\{I_j\} = p_j$.

In a recent paper (Praskova, 1985) the weak convergence of S_{N_n} to a Poisson r.v. as $N, n \rightarrow \infty$, was discussed in some detail. In this note, we extend and generalize the results of Praskova (1985). First, we show that the Poisson convergence still holds when $\{Y_j\}_1^N$ are non-negative independent integer valued random variables (r.v.'s), independent of $\{I_j\}_1^N$, such that $E\{Y_j\}^k < \infty, k = 1, 2$. The randomness of Y_j covers cases such as multistage sampling where Y_j is the value corresponding to the j th primary stage unit. Second, under mild regularity assumptions on $\{Y_j\}_1^N$, we also investigate the weak convergence of the standardized sample sum to a normal distribution.

2. Poisson convergence of the sample sum

Set $P_{N_n} = P\{\sum_1^N I_j = n\}$, $f_j(t) = E\{e^{itY_j}\}$ and

$$\varphi_{N_n}(t) = E\left(\exp\left\{it\sum_1^N Y_j I_j\right\} \middle| \sum_1^N I_j = n\right).$$

From a simple argument (see e.g. Holst, 1979, Theorem 1) we have:

$$\varphi_{N_n}(t) = (2\pi P_{N_n})^{-1} \int_{-\pi}^{\pi} e^{-isn} \left[\prod_{j=1}^N E\left(\exp\{i(tY_j + s)I_j\}\right) \right] ds. \tag{2.1}$$

* Permanent Address: Department of Mathematics, Eötvös University, Budapest, Hungary.

Let Z_N be a Poisson r.v. with $E\{Z_N\} = \sum_1^N E(Y_j)p_j$, then its characteristic function $g_N(t)$ is

$$g_N(t) = \exp\left\{ (e^{it} - 1) \sum_1^N p_j E(Y_j) \right\}. \tag{2.2}$$

Our aim is to evaluate $|\varphi_{N_n}(t) - g_N(t)|$, which we do in Proposition 2.1. Throughout this paper, we assume that $p_j < 0.5$. Set $q_j = 1 - p_j$, then

$$E\left\{ \exp\{i(tY_j + s)I_j\} \right\} = (q_j + p_j e^{is} f_j(t)). \tag{2.3}$$

This and (2.1) yield

$$\varphi_{N_n}(t) = (2\pi P_{N_n})^{-1} \int_{-\pi}^{\pi} e^{-isn} \left[\prod_1^N (q_j + p_j e^{is}) (q_j(s) + p_j(s) f_j(t)) \right] ds \tag{2.4}$$

where $q_j(s) = q_j / (q_j + p_j e^{is})$ and $p_j(s) = 1 - q_j(s)$. Since

$$1 = (2\pi P_{N_n})^{-1} \int_{-\pi}^{\pi} e^{-isn} E\left\{ \exp\left\{ is \sum_1^N I_j \right\} \right\} ds$$

it follows, using a well known identity for products, that

$$\begin{aligned} \varphi_{N_n}(t) - g_N(t) &= (2\pi P_{N_n})^{-1} \int_{-\pi}^{\pi} e^{-isn} \prod_1^N (q_j + p_j e^{is}) \\ &\quad \times \left\{ \prod_1^N (q_j(s) + p_j(s) f_j(t)) - \prod_1^N e^{p_j(e^{it}-1)E(Y_j)} \right\} ds \\ &= (2\pi P_{N_n})^{-1} \int_{-\pi}^{\pi} e^{-isn} \prod_1^N (q_j + p_j e^{is}) \\ &\quad \times \left\{ \sum_{j=1}^N \left[\prod_{l=1}^{j-1} (q_l(s) + p_l(s) f_l(t)) \right] \left[\prod_{k=j+1}^N e^{p_k(e^{it}-1)E(Y_k)} \right] \right. \\ &\quad \left. \times [q_j(s) + p_j(s) f_j(t) - e^{p_j(e^{it}-1)E(Y_j)}] \right\} ds. \tag{2.5} \end{aligned}$$

To evaluate the absolute value of the difference in (2.5), the following lemma is needed.

Lemma 2.1.

$$\begin{aligned} &|q_j(s) + p_j(s) f_j(t) - e^{p_j(e^{it}-1)E(Y_j)}| \\ &\leq |e^{it} - 1|^2 \left[\frac{p_j E Y_j (Y_j - 1)}{2(q_j - p_j)} + \frac{2 p_j q_j |\sin \frac{1}{2}s| E(Y_j)}{(q_j - p_j) |e^{it} - 1|} + \frac{1}{2} (p_j E(Y_j))^2 \right]. \tag{2.6} \end{aligned}$$

Proof. First, write the left hand side of (2.6) as follows:

$$\begin{aligned} & \left| p_j(s) \left[f_j(t) - 1 - E(Y_j)(e^{it} - 1) \right] + (p_j(s) - p_j)(e^{it} - 1)E(Y_j) \right. \\ & \quad \left. - \left(e^{p_j(e^{it}-1)E(Y_j)} - 1 - p_j(e^{it} - 1)E(Y_j) \right) \right|. \end{aligned} \tag{2.7}$$

Clearly,

$$\begin{aligned} |p_j(s)| &\leq \frac{p_j}{q_j - p_j}, & |p_j(s) - p_j| &\leq \frac{2p_jq_j}{q_j - p_j} |\sin \frac{1}{2}s|, \\ |f_j(t) - 1 - E(Y_j)(e^{it} - 1)| &= \left| \sum_{k=0}^{\infty} e^{ikt} P\{Y_j = k\} - 1 - (e^{it} - 1) \sum_{k=0}^{\infty} kP\{Y_j = k\} \right| \\ &\leq |e^{it} - 1| \sum_{k=0}^{\infty} \left| \left[\frac{e^{ikt} - 1}{e^{it} - 1} - k \right] \right| P\{Y_j = k\} \\ &\leq \frac{1}{2} |e^{it} - 1|^2 EY_j(Y_j - 1). \end{aligned}$$

Finally,

$$\left| e^{p_j(e^{it}-1)E(Y_j)} - 1 - p_j(e^{it} - 1)E(Y_j) \right| \leq \frac{1}{2} |e^{it} - 1|^2 (p_j E(Y_j))^2.$$

This proves the lemma. \square

Set

$$\begin{aligned} A &= \sum_1^N \frac{p_j E(Y_j)}{q_j - p_j}, & B &= \sum_1^N \left[\frac{p_j E Y_j (Y_j - 1)}{q_j - p_j} + (p_j E(Y_j))^2 \right], \\ C &= \sum_1^N \frac{p_j q_j}{q_j - p_j} E\{Y_j\}, & d &= \sum_1^N p_j q_j. \end{aligned}$$

Then we have:

Proposition 2.1.

$$|\varphi_{N_n}(t) - g_N(t)| \leq \alpha e^{|\epsilon^{it}-1|A} \cdot \{ \beta |e^{it} - 1| B + \gamma |e^{it} - 1| d^{-1/2} C \} \tag{2.8}$$

where α, β and γ are positive constants.

Proof. First, we have the following three inequalities:

- (i) $|q_l(s) + p_l(s)f_l(t)| = |1 + (f_l(t) - 1)p_l(s)| \leq 1 + |e^{it} - 1| p_l(s) E(Y_l) \leq \exp\{ p_l E(Y_l) |e^{it} - 1| / (q_l - p_l) \}.$
- (ii) $|e^{p_k(e^{it}-1)E(Y_k)}| \leq e^{p_k |e^{it}-1| E(Y_k) / (q_k - p_k)}.$
- (iii) $|q_j + p_j e^{is}| \leq e^{-2p_j q_j \sin^2(s/2)}.$

This and (2.5) yield:

$$\begin{aligned} & |\varphi_{N_n}(t) - g_N(t)| \\ & \leq (2\pi P_{N_n})^{-1} e^{|\epsilon^{it}-1|A} \int_{-\pi}^{\pi} e^{-2d \sin^2(s/2)} \left[\frac{1}{2} B |e^{it} - 1|^2 + C |e^{it} - 1| |\sin \frac{1}{2}s| \right] ds. \end{aligned}$$

Since

$$\int_0^\pi e^{-2d \sin^2(s/2)} ds \leq \left[\frac{1}{2}\pi\right]^{3/2} d^{-1/2}, \quad \int_0^\pi e^{-2d \sin^2(s/2)} \sin \frac{1}{2}s ds \leq (1 - e^{-2d})/d,$$

and $P_{N_n} \geq k \cdot d^{1/2}$, where $k > 0$ is a constant, the proposition follows. \square

An important consequence of the proposition is the Poisson convergence of the sample sum, stated formally in Corollary 2.1. Set $(\alpha = \alpha(N))$ is assumed to be a bounded sequence

$$\alpha^{-1} = \min_{1 \leq j \leq N} (q_j - p_j).$$

Then clearly

$$\begin{aligned} A &\leq \alpha \sum_1^N p_j E\{Y_j\}, \\ B &\leq \alpha \sum_1^N p_j E Y_j (Y_j - 1) + \max_{1 \leq k \leq N} (p_k E(Y_k)) \sum_1^N p_j E(Y_j), \\ C &\leq \alpha \sum_1^N p_j E\{Y_j\}. \end{aligned}$$

In addition, since $\sum_1^N p_j = n$ and $p_j < 0.5$ it follows that $d \geq 2n$.

Corollary 2.1. *If we set $I_j = I_{N_j}$, $p_j = p_{N_j}$, and assume that as $N \rightarrow \infty$, $n \rightarrow \infty$,*

$$\max_{1 \leq j \leq N} p_{N_j} E(Y_j) \rightarrow 0, \quad \sum_1^N p_{N_j} E Y_j (Y_j - 1) \rightarrow 0$$

and

$$\sum_1^N p_{N_j} E(Y_j) \rightarrow \lambda,$$

then $\varphi_{N_n}(t) \rightarrow g(t)$, where $g(t) = e^{(e^t - 1)\lambda}$. This implies that

$$\mathcal{L} \left\{ \sum_1^N Y_j I_{N_j} \mid \sum_1^N I_{N_j} = n \right\} \rightarrow \mathcal{L}(Z)$$

where Z is a Poisson r.v. with $E\{Z\} = \lambda$. \square

3. Convergence to a normal distribution

The result of this section can be formulated as:

Proposition 3.1. *Let $\{Y_j\}_1^\infty$ be independent real-valued r.v.'s such that $E\{Y_j\} = 0$ and $\text{Var}\{Y_j\} = 1$, $j = 1, 2, \dots$, independent of $\{I_j\}_1^\infty$. Then*

$$S_{N_n} / \sqrt{n} \xrightarrow{\mathcal{L}} Z$$

where $Z \sim N(0,1)$, or equivalently $\Psi_{N_n}(t) \rightarrow \phi(t)$, uniformly where

$$\Psi_{N_n}(t) = E \left(\exp \left\{ \frac{it}{\sqrt{n}} \sum_1^N Y_j I_j \middle| \sum_1^N I_j = n \right\} \right) \quad \text{and} \quad \phi(t) = e^{-t^2/2}.$$

Proof. The method of proof is the one used in the previous proposition. Write

$$\Psi_{N_n}(t) = (2\pi P_{N_n})^{-1} \int_{-\pi}^{\pi} e^{-isn} \prod_1^N \left(q_j + p_j e^{is} f_j \left(\frac{t}{\sqrt{n}} \right) \right) ds.$$

Then, we have

$$\begin{aligned} \Psi_{N_n}(t) - \phi(t) &= (2\pi P_{N_n})^{-1} \int_{-\pi}^{\pi} e^{-isn} \left[\prod_1^N (q_j + p_j e^{is}) \right] \\ &\quad \times \left[\prod_1^N \left(q_j(s) + p_j(s) f_j \left(\frac{t}{\sqrt{n}} \right) \right) - \prod_1^N e^{-(t^2/2n)p_j} \right] ds \\ &= (2\pi P_{N_n})^{-1} \int_{-\pi}^{\pi} e^{-isn} \left[\prod_1^N (q_j + p_j e^{is}) \right] \\ &\quad \times \left\{ \sum_{j=1}^N \left[\prod_{k=1}^{j-1} \left(q_k(s) + p_k(s) f_k \left(\frac{t}{\sqrt{n}} \right) \right) \right] \left[\prod_{k=j+1}^N e^{-(t^2/2n)p_k} \right] \right. \\ &\quad \left. \times \left(q_j(s) + p_j(s) f_j \left(\frac{t}{\sqrt{n}} \right) - e^{-(t^2/2n)p_j} \right) \right\} ds. \end{aligned}$$

But

$$\begin{aligned} &\left| q_j(s) + p_j(s) f_j \left(\frac{t}{\sqrt{n}} \right) - e^{-(t^2/2n)p_j} \right| \\ &= \left| p_j(s) \left(f_j \left(\frac{t}{\sqrt{n}} \right) - 1 \right) + (1 - e^{-(t^2/2n)p_j}) \right| \\ &= \left| p_j(s) \left(f_j \left(\frac{t}{\sqrt{n}} \right) \right) + (p_j - p_j(s)) \frac{t^2}{2n} + 1 - \frac{t^2}{2n} p_j - e^{-(t^2/2n)p_j} \right| \\ &\leq \frac{p_j}{q_j - p_j} o\left(\frac{1}{n}\right) + \frac{2p_j q_j |\sin \frac{1}{2}s|}{q_j - p_j} \frac{t^2}{2n} + o\left(\frac{1}{n}\right) \frac{t^2}{2} p_j, \\ &\left| q_k(s) + p_k(s) f_k \left(\frac{t}{\sqrt{n}} \right) \right| = \left| 1 + p_k(s) \left(f_k \left(\frac{t}{\sqrt{n}} \right) - 1 \right) \right| \\ &\leq \exp \left\{ |p_k(s)| \cdot \left| f_k \left(\frac{t}{\sqrt{n}} \right) - 1 \right| \right\} \leq \exp \left\{ \frac{p_k}{q_k - p_k} o\left(\frac{t^2}{2n}\right) \right\}, \end{aligned}$$

and

$$\left| q_j + p_j e^{is} \right| \leq \exp \left\{ 2p_j q_j \sin^2 \frac{1}{2}s \right\}.$$

From this, we obtain that

$$\begin{aligned}
 |\Psi_{N,n}(t) - \phi(t)| &\leq (2\pi P_{N,n})^{-1} \int_{-\pi}^{\pi} e^{-(\sum_1^N q_j p_j) \sin^2(s/2) - (t^2/2n) \sum_1^N (p_k(q_k - p_k)) + o(1)} \\
 &\quad \times \left[o(1) + \frac{t^2}{2n} \left(\sum_1^N q_j p_j \right) \sin \frac{1}{2}s + o(1) \right] ds \\
 &\leq d^{-1/2} \left[\frac{o(1)}{\sqrt{d}} + \frac{t^2}{2n} \right] = o(1) + \frac{c}{\sqrt{d}} \rightarrow 0
 \end{aligned}$$

which proves the assertion. \square

References

- Hajek, J. (1981), *Sampling from a Finite Population* (Dekker, New York).
- Holst, L. (1979), Two conditional limit theorems with applications, *Ann. Statist.* **7**, 551–557.
- Praskova, Z. (1985), The convergence to the Poisson distribution in rejective sampling from a finite population, in: W. Grosman, J. Mogyorodi, I. Vincze and W. Wertz, eds., *Proc. 5th Panonian Symp. on Math. Statist., Visegrad Hungary, 1985*.