

SIMPLE APPROXIMATIONS TO THE VARIATION NORM  
BETWEEN SAMPLING WITH AND WITHOUT REPLACEMENT

Chai Ho C. Wang, U.S. Department of Justice

The variation norm between sampling with and without replacement is given by Freedman [1]. Drawing a sample of  $k$  from a population of  $N$  ( $0 < k \leq N$ ), the variation norm is given by:

$$d = \|P - Q\| = 1 - (N_k / N^k)$$

where  $P$ ,  $Q$  are the induced probabilities for sampling with and without replacement respectively, and

$$N_k = N(N-1) \cdots (N-k+1)$$

Freedman's bounds for this norm are:

$$l_1 = 1 - \exp[-k(k-1)/2N] \\ l_1 < d < k(k-1)/2N = u_1$$

While  $l_1$  is a fairly accurate approximation for  $d$ ,  $u_1$ , as the first term of the Taylor series for  $l_1$ , is less adequate in representing  $d$  for certain combinations of  $N$  and  $k$ .

A simple approximation for  $d$  can be obtained via Stirling's formula for  $n!$ :

$$n! \approx s_n = 2\pi n(n/e)^n$$

We have

$$\frac{N_k}{N^k} = \frac{N!}{(N-k)!N^k} \\ \approx e^{-k} \left(\frac{N}{N-k}\right)^{N-k+1/2}$$

and

$$l_3 = 1 - e^{-k} \left(\frac{N}{N-k}\right)^{N-k+1/2} \\ l_3 < d, N \neq k$$

As will be justified later this formula provides accurate approximation for  $d$ . However, due to the involvement of large exponents, computational inconvenience and undesirable rounding error may arise. Taking an alternative approach, let us write, say, for  $k$  odd,

$$N(k) = \frac{N_k}{N^k} = \prod_{j=1}^{k-1} (1-j/N) \\ = \prod_{j=1}^{(k-1)/2} \left(1 - \frac{j}{N}\right) \left(1 - \frac{k-j}{N}\right) \\ = \prod_{j=1}^{(k-1)/2} \left(1 - \frac{k}{N} + \frac{j(k-j)}{N^2}\right) \quad (I)$$

Observing that

$$\min_{1 \leq j \leq k-1} \frac{j(k-j)}{N^2} = \frac{k-1}{N^2},$$

$$\max_{1 \leq j \leq k-1} \frac{j(k-j)}{N^2} = \frac{k^2-1}{4N^2}$$

we have

$$\left(1 - \frac{k}{N} + \frac{k^2-1}{N^2}\right)^{(k-1)/2} < N(k) \\ \left(1 - \frac{k}{N} + \frac{k^2-1}{4N^2}\right)^{(k-1)/2}$$

Further simplifications lead to the following two bounds for  $d$ :

$$l_2 = 1 - \left(1 - \frac{k}{2N}\right)^k < d$$

and

$$u_2 = 1 - \left(1 - \frac{k}{N}\right)^{k/2} > d$$

where  $k =$  the greatest odd integer less than or equal to  $k$ , and  $k =$  the greatest even integer less than or equal to  $k$ .

Since

$$\sum_{j=1}^{(k-1)/2} \frac{j(k-j)}{N^2} = \frac{k(k^2-1)}{12N^2},$$

sharper approximations can be obtained by expanding (I). We have

$$u_3 = 1 - (1-k/N)^{(k-1)/2} \left[1 - \frac{k}{N} + \frac{k(k^2-1)}{12N^2}\right] > d$$

For practical computation purpose where  $N$  is large,  $u_3$  is sufficiently accurate that no measurable precision will be gained by calculating the exact solution. To pursue this argument one step further, we note that the leading order term discarded when deriving  $u_3$ , again for  $k$  odd, is bounded above by

$$\hat{u}_3 = (1-k/N)^{(k-1)/2} \prod_{i<j}^{(k-5)/2} \frac{i(k-i) \cdot j(k-j)}{N^4} \\ < i/2(1-k/N)^{(k-5)/2} \left[ \sum_{j=1}^{(k-1)/2} \frac{j(k-j)}{N^2} \right]^2 \\ < (1-k/N)^{(k-5)/2} \frac{k^6}{288N^4}$$

This upper bound for  $\hat{u}_3$  can be maximized at  $k \approx \sqrt{12N}$ . It follows that

$$\hat{u}_3 < 6g(N)/N$$

where

$$g(N) = (1 - \sqrt{12/N}) (\sqrt{12N} - 5)/2$$

$g(N)$  is an increasing function of  $N$ . It increases very slowly for  $N$  up to, say,  $10^{20}$ . Consequently,

$$\hat{u} < \frac{1}{N}$$

for all  $k$ , and for almost all  $N$ .

Let the relative error for approximating  $d$  by  $d'$  be defined by  $R = (d' - d)/d'$ . The difference between  $R$  and the usual relative error  $r = (d' - d)/d$  is small if  $R$  and  $r$  are small, because  $|r| < |R| \leq |r-R| = |rR|$ . The analysis of the truncation error for  $u_3$  can be extended to the following proposition:

Proposition. The relative error for calculating  $N_k/N^k$  while using the (approximation) formulas associated with  $l_1, u_1, l_2, u_2, l_3$ , and  $u_3$  are  $O[k^3/N(N-k)]$ ,  $O[k^4/(N(2N-k^2))]$ ,  $O[k^3/N(N-k)]$ ,  $O[k^3/N(N-k)]$ ,  $O[k/N(N-k)]$ , and  $O[k^6/N^2(N-k)^2]$  respectively.

Proof: Since

$$\frac{\hat{u}_3}{u_3} < \frac{\hat{u}_3}{\frac{(k-1)/2}{(1-k/N)}} < \frac{k^6}{288N^2(N-k)^2}$$

the proposition for the case  $u_3$  has been established. Similar leading-order-term argument is sufficient for proving all other cases except the case  $l_3$ . For this case we note that

$$\ln(n!) - \ln s_n = \frac{1}{12n} - \frac{1}{360n^2} + \dots < \frac{1}{12n}$$

then

$$\begin{aligned} N!/s_N \\ 0 < R = 1 - \frac{N!}{(N-k)!/s_{N-k}} \\ &< 1 - \exp\left(\frac{1}{12N} - \frac{1}{12(N-k)}\right) \\ &< \frac{k}{12N(N-k)} \end{aligned}$$

as was to be shown.

From this proposition, it is apparent that both  $l_3$  and  $u_3$  provide accurate approximation for  $d$ . Although as indicated in the proposition,  $l_3$  should be favored whenever  $k^5 > 24N(N-K)$ , only a rigorous analysis of rounding errors, which is outside the scope of this paper, can decide whether and when  $l_3$  is computationally more useful than  $u_3$ .

<sup>3</sup> Table 1 displays numerical results computed with 16 (decimal) digit accuracy. The "exact" solution is computed by brute force using

$$\hat{d} = 1 - \prod_{j=1}^{k-1} \left(1 - \frac{j}{N}\right)$$

For a fixed  $N$ , as  $k$  increases, all approximations except  $u_1$  rapidly approach unity.

Relative computing error

$$\hat{R} = \hat{d} - d' = R + \text{rounding error}$$

for these approximations are found in Table 2. In certain cases, the rounding error becomes dominate

over the relative error. For computation purpose  $u_3$  appears to be a better approximation.

#### Reference

1. Freedman, D. (1977), "A Remark on the Difference Between Sampling With and Without Replacement" *Journal of the American Statistical Association*, Volume 72, 681.

Table 1

Approximations of  $d = 1 - N_k/N^k$   
 $[a(J) = a \times 10^J]$

Table 2  
Relative Computing Error