If the observations are independent and identically distributed, the variance of ft. is given by

$$v(\ ) \ n^{-1}0^{,2} \tag{13}$$

which can be estimated by replacing $c_{r^2}$ by its estimate from (12). Consider now what happens when the x's are no longer i.i.d., but belong to clusters, and that within each cluster

$$E(x,\ IL)(x_i\ 1^4) = Pa^2 \tag{14}$$

for some quantity *p,* while for two observations in different clusters, we retain the assumption of independence. Then, as shown by Kish (1965), and as may be readily confirmed, (13) must be replaced by

$$V(4) = n\ cr^2\ d \tag{15}$$

where *d* is the Kish design effect, or "deff", defined by

$$d\ 1 + (ii,\ —1)p\ . \tag{16}$$

The quantity $i_t$ is the number of households in each cluster when the clusters are all the same size; more generally it is the weighted average of cluster sizes, where the weights are the cluster sizes themselves, i.e. n $E/2,2$ for individual cluster sizes n,. An estimate of *p* can be obtained from the "intracluster correlation coefficient"

$$\frac{E,\qquad (x_{,,,}\ 1.1)(x_{j,}}{6^{-2}\ E,\ nc(n,\ —1)}\qquad \bullet \tag{17}$$

A number of points should be noted. In the presence of positive intracluster correlations, the number of "effective" observations is smaller than the sample size. In the extreme case, when $_p$ is unity, *d is* the cluster size, and the effective sample size is the number of clusters, not the number of observations. Even when *p* is 0.5, a high but not unusual figure, the usual formula for the standard error of a mean is optimistic by a factor of 2.34 (the square root of 5), a correction that could make a substantial difference to the conclusions being drawn. Second, although I have illustrated using clusters, the same analysis might be useful within strata, or regions, or sectors, or any other partition of the sample for which there is reason to believe that the observations within each partition are correlated. When the partition is large, *p is* likely to be small, but the size of "deff" depends on the product, and might still be large.