

CHAPTER 23

Estimation of Distribution Functions and Quantiles

23.1 INTRODUCTION

In earlier sections, we have considered various methods of estimating the finite population totals, means, and ratios of two variables. Most of these methods, when extended to the estimation of distribution function, yield unsatisfactory results because these estimators may not satisfy the basic properties of the distribution functions. Estimation of distribution functions from the survey data is often a part of the objectives of a survey. In particular, it may be of interest to planners who would want to know the proportion of people living below the poverty line, unemployed, or money spent on education. Furthermore, it is well known that the median is considered to be a more appropriate measure of location than the mean of a skewed distribution such as income. It is also important to estimate the income inequality through the Gini coefficient, which is a function of a distribution function. In this section, different methods of estimation of distribution function have been considered. Estimation of quantiles and medians will also be obtained from the inversion of the estimates of distribution functions.

23.2 ESTIMATION OF DISTRIBUTION FUNCTIONS

The distribution function of a random variable X is defined as $F(x) = P(X \leq x)$. The distribution function $F(x)$ has the following properties:

- (i) $F(x) \geq 0$, (ii) nondecreasing, (iii) $F(-\infty) = 0$, (iv) $F(+\infty) = 1$, and (v) $F(x)$ is right continuous

(23.2.1)

For a finite population of N units, we define the distribution function $F(t)$ of the study variable y , which is the proportion of the units of the population U whose values are less than or equal to t , i.e.,

$$F(t) = \frac{1}{N} \sum_{i=1}^N I(y_i \leq t) \quad (23.2.2)$$

where $I(y_i \leq t) = \begin{cases} 1 & \text{if } y_i \leq t \\ 0 & \text{if } y_i > t \end{cases}$ and y_i is the value of the study variable y for the i th unit of the population; $i = 1, \dots, N$.

23.2.1 Design-Based Estimation

Let a sample s of size n be selected from a population U using a sampling design p with $\pi_i(>0)$ and $\pi_{ij}(>0)$ as the inclusion probabilities for the i th, and i th and j th unit ($i \neq j$), respectively. In case the population size N is known, an unbiased estimator of $F(t)$ is given by

$$\hat{F}(t) = \frac{1}{N} \sum_{i \in s} \frac{I(y_i \leq t)}{\pi_i} \quad (23.2.3)$$

For simple random sampling without replacement (SRSWOR) sampling $\pi_i = n/N$ and $\hat{F}(t)$ reduces to the sample empirical distribution function

$$s(t) = \frac{r(t)}{n} \quad (23.2.4)$$

where $r(t)$ is number of y_i 's in the sample less than or equal to t .

The estimator $s(t)$ is admissible and also minimax under certain loss function (Dorfman, 2009). KuK (1988) proposed following alternative estimators based on complementary proportion

$$\hat{F}_c(t) = 1 - \frac{1}{N} \sum_{i \in s} \frac{1 - I(y_i \leq t)}{\pi_i} \quad (23.2.5)$$

and

$$\hat{F}_w(t) = w\hat{F}(t) + (1 - w)\hat{F}_c(t) \quad (23.2.6)$$

where w is a suitably chosen weight.

It should be noted that none of the estimators $\hat{F}(t)$, $\hat{F}_c(t)$, and $\hat{F}_w(t)$ mentioned above possess all the properties of the distribution function stated in Eq. (23.2.1). In fact $\hat{F}(t)$ may exceed unity and $\hat{F}(\infty) = \frac{1}{N} \sum_{i \in s} \frac{1}{\pi_i}$ may

not equal to 1. Details on discussions of the performances of the estimators have been given by Dorfman (2009). To overcome this difficulty, the following alternative Hájek type estimator has been proposed.

$$\hat{F}_h(t) = \sum_{i \in s} \frac{I(y_i \leq t)}{\pi_i} \bigg/ \sum_{i \in s} \frac{1}{\pi_i} \quad (23.2.7)$$

Though the estimator (23.2.7) is not design unbiased for $F(t)$, it is design consistent and approximately unbiased. It satisfies all the properties of the distribution function stated in Eq. (23.2.1).

23.2.2 Design-Based Estimators Using Auxiliary Information

Rao et al. (1990) considered design-based ratio and difference estimators of the population distribution function $F(t)$ when the values of the auxiliary variable x_i 's are known for $i = 1, \dots, N$. The conventional estimator for the population ratio $R = Y/X$ is

$$\hat{R} = \left(\sum_{i \in s} \frac{y_i}{\pi_i} \right) / \left(\sum_{i \in s} \frac{x_i}{\pi_i} \right) \quad (23.2.8)$$

where $X = \sum_{i=1}^N x_i$ and $Y = \sum_{i=1}^N y_i$.

Treating $\hat{y}_i = \hat{R}x_i$ as an estimator of y_i , Rao et al. (1990) proposed the following ratio estimator of $F(t)$ as

$$\hat{F}_R(t) = \frac{1}{N} \hat{\lambda} \times \left(\sum_{i \in U} I(\hat{y}_i \leq t) \right) \quad (23.2.9)$$

where

$$\hat{y}_i = \hat{R}x_i \quad \text{and} \quad \hat{\lambda} = \left(\frac{\sum_{i \in s} \frac{I(y_i \leq t)}{\pi_i}}{\sum_{i \in s} \frac{I(\hat{y}_i \leq t)}{\pi_i}} \right) \quad (23.2.10)$$

The proposed ratio estimator $\hat{F}_R(t)$ is design consistent and exactly equal to the population distribution function $F(t)$ if y_i is proportional to x_i for every $i \in U$. The proposed estimator is expected to gain in efficiency compared to the conventional estimator $\hat{F}(t)$ if y_i is approximately proportional to x_i . But the estimator $\hat{F}_R(t)$ suffers from a drawback that it may not possess the desirable properties of the distribution function given in Eq. (23.2.1): it may have a value outside the interval $[0, 1]$. For a large sample size n , an approximate expression of variance of $\hat{F}_R(t)$ is given by Rao et al. (1990) as

$$\begin{aligned} \text{Var}[\hat{F}_R(t)] &\cong V \left(\frac{1}{N} \sum_{i \in s} \frac{Q_i(t)}{\pi_i} \right) \\ &= \frac{1}{2N^2} \sum_{i \neq j} \sum_{j \in U} (\pi_i \pi_j - \pi_{ij}) \left(\frac{Q_i}{\pi_i} - \frac{Q_j}{\pi_j} \right)^2 \end{aligned} \quad (23.2.11)$$

where

$$Q_i = I(\gamma_i \leq t) - \lambda I(Rx_i \leq t) \text{ and } \lambda = \frac{\sum_{i \in U} I(\gamma_i \leq t)}{\sum_{i \in U} I(Rx_i \leq t)}. \quad (23.2.12)$$

An approximate unbiased estimator of $V_p[\hat{F}_R(t)]$ is

$$\hat{V}[\hat{F}_R(t)] = \frac{1}{2N^2} \sum_{i \neq j} \sum_{j \in s} \frac{(\pi_i \pi_j - \pi_{ij})}{\pi_{ij}} \left(\frac{\hat{Q}_i}{\pi_i} - \frac{\hat{Q}_j}{\pi_j} \right)^2 \quad (23.2.13)$$

where $\hat{Q}_i = I(\gamma_i \leq t) - \hat{\lambda} I(\hat{\gamma}_i \leq t)$ and $\hat{\lambda}$ is given in Eq. (23.2.10).

Rao et al. (1990) also proposed a difference estimator of $F(t)$ as

$$\hat{F}_D(t) = \frac{1}{N} \left\{ \left(\sum_{i \in s} \frac{I(\gamma_i \leq t)}{\pi_i} \right) - \left(\sum_{i \in s} \frac{I(\hat{\gamma}_i \leq t)}{\pi_i} - \sum_{i \in U} I(\hat{\gamma}_i \leq t) \right) \right\} \quad (23.2.14)$$

An approximate expression of the variance of $\hat{F}_D(t)$ and its unbiased estimators are respectively given by

$$\begin{aligned} Var[\hat{F}_D(t)] &\cong V_p \left(\frac{1}{N} \sum_{i \in s} \frac{G_i(t)}{\pi_i} \right) \\ &= \frac{1}{2N^2} \sum_{i \neq j} \sum_{j \in U} (\pi_i \pi_j - \pi_{ij}) \left(\frac{G_i(t)}{\pi_i} - \frac{G_j(t)}{\pi_j} \right)^2 \end{aligned} \quad (23.2.15)$$

and

$$\hat{V}[\hat{F}_D(t)] = \frac{1}{2N^2} \sum_{i \neq j} \sum_{j \in s} \frac{(\pi_i \pi_j - \pi_{ij})}{\pi_{ij}} \left(\frac{\hat{G}_i(t)}{\pi_i} - \frac{\hat{G}_j(t)}{\pi_j} \right)^2 \quad (23.2.16)$$

where $G_i(t) = I(\gamma_i \leq t) - I(Rx_i \leq t)$ and $\hat{G}_i(t) = I(\gamma_i \leq t) - I(\hat{R}x_i \leq t)$.

23.2.3 Model-Based Estimators

Let us suppose that the study variable y is related to the auxiliary variable through the superpopulation model

$$y_i = \mu(x_i) + \sqrt{v(x_i)} \epsilon_i \quad (23.2.17)$$

where $\mu(x_i)$ is a function of x_i but involves unknown model parameters, $\nu(x_i)$ is known and positive, and ϵ_i 's are independently and identically distributed random variables with $E_m(\epsilon_i) = 0$ and $V_m(\epsilon_i) = \sigma^2(>0)$. Here E_m and V_m denote expectation and variance operators, respectively, with respect to the model (23.2.17). Chambers and Dunstan (1986) and Rao et al. (1990) used the model (23.2.17) with $\mu(x_i) = \beta x_i$, where β is an unknown parameter.

Let $\hat{\mu}(x_i)$ be a model-unbiased estimator of $\mu(x_i)$, i.e., $E_m[\hat{\mu}(x_i)] = \mu(x_i)$. Then $\hat{r}_i(t) = \frac{t - \hat{\mu}(x_i)}{\sqrt{\nu(x_i)}}$ becomes model-unbiased estimator of $r_i(t) = \frac{t - \mu(x_i)}{\sqrt{\nu(x_i)}}$.

Let $A_i(t)$ be the model expectation of $I(y_i \leq t)$ in the sense that

$$A_i(t) = E_m\{I(y_i \leq t)\} = P(y_i \leq t) = P(\epsilon_i \leq r_i(t)) \quad (23.2.18)$$

An approximate model-unbiased estimator of $A_i(t)$ is given by

$$\hat{A}_i(t) = \frac{1}{n} \sum_{j \in s} I\{\hat{r}_j(y_j) \leq \hat{r}_i(t)\} \quad (23.2.19)$$

where $\hat{r}_j(y_j) = \frac{y_j - \hat{\mu}(x_j)}{\sqrt{\nu(x_j)}}$.

Under the model (23.2.17), Johnson (2003) proposed the following model-based estimator of $F(t)$ as

$$\hat{F}_J(t) = \frac{1}{N} \left(\sum_{i \in s} I(y_i \leq t) + \sum_{i \in U-s} \hat{A}_i(t) \right) \quad (23.2.20)$$

where $U - s$ denotes the set of nonsampled units.

Chambers and Dunstan (1986) considered the estimator (23.2.20) when $\mu(x_i) = \beta x_i$ and the model parameter β was estimated by the weighted least squares method as

$$\hat{\beta} = \left(\sum_{i \in s} \frac{y_i x_i}{\nu(x_i)} \right) / \left(\sum_{i \in s} \frac{x_i^2}{\nu(x_i)} \right) \quad (23.2.21)$$

Substituting $\hat{\mu}(x_i) = \hat{\beta} x_i$ in the expression (23.2.20), Chambers and Dunstan (1986) estimator is obtained as

$$\hat{F}_{cd}(t) = \frac{1}{N} \left[\sum_{i \in s} I(y_i \leq t) + \sum_{i \in U-s} \left\{ \frac{1}{n} \sum_{j \in s} I\left(\frac{y_j - \hat{\beta} x_j}{\sqrt{\nu(x_j)}} \leq \frac{t - \hat{\beta} x_i}{\sqrt{\nu(x_i)}} \right) \right\} \right] \quad (23.2.22)$$

The estimator \hat{F}_{cd} satisfies all the properties of a distribution function. Although the estimator $\hat{F}_{cd}(t)$ is independent of the sampling design, it is asymptotically model unbiased and has the property that $\hat{F}_{cd}(t)$ equals to $F(t)$ when $y_i \propto x_i \forall i \in U$. Since the estimator \hat{F}_{cd} depends on the assumed model, it is highly efficient if the model (23.2.17) holds but if the model is incorrect, the estimator becomes biased and can perform much worse even than the naive estimator $s(t)$.

23.2.4 Model-Assisted Estimators

Rao et al. (1990) proposed the following model assisted difference estimators:

$$\hat{F}^*(t) = \frac{1}{N} \left[\sum_{i \in s} \frac{I(y_i \leq t)}{\pi_i} - \left(\sum_{i \in s} \frac{B_i(t)}{\pi_i} - \sum_{i \in U} B_i(t) \right) \right] \quad (23.2.23)$$

with

$$B_i(t) = \frac{1}{N} \sum_{j \in U} I \left(\frac{y_j - \beta x_j}{\sqrt{v(x_j)}} \leq \frac{t - \beta x_i}{\sqrt{v(x_i)}} \right) \quad (23.2.24)$$

The estimator $\hat{F}^*(t)$ does not possess all the properties of distribution function. It is both design unbiased and asymptotically model unbiased for $F(t)$. It is calibrated in the sense that $\hat{F}^*(t)$ reduces $F(t)$ if y_i is exactly proportional to x_i for $\forall i \in U$. The main demerit of this estimator is that it cannot be used in practice as $B_i(t)$'s are unknown. To overcome this difficulty, Rao et al. (1990) proposed the following alternative estimator:

$$\hat{F}_{rkm}(t) = \frac{1}{N} \left[\sum_{i \in s} \frac{I(y_i \leq t)}{\pi_i} - \left(\sum_{i \in s} \frac{\hat{B}_{ic}(t)}{\pi_i} - \sum_{i \in U} \hat{B}_i(t) \right) \right] \quad (23.2.25)$$

where

$$\begin{aligned} \hat{B}_i(t) &= \left(\sum_{j \in s} \frac{1}{\pi_j} \right)^{-1} \left[\sum_{j \in s} \frac{1}{\pi_j} I \left(\frac{y_j - \hat{\beta} x_j}{\sqrt{v(x_j)}} \leq \frac{t - \hat{\beta} x_i}{\sqrt{v(x_i)}} \right) \right] \text{ and} \\ \hat{B}_{ic}(t) &= \left(\sum_{j \in s} \frac{\pi_j}{\pi_{ij}} \right)^{-1} \left[\sum_{j \in s} \frac{\pi_j}{\pi_{ij}} I \left(\frac{y_j - \hat{\beta} x_j}{\sqrt{v(x_j)}} \leq \frac{t - \hat{\beta} x_i}{\sqrt{v(x_i)}} \right) \right] \end{aligned} \quad (23.2.26)$$

The estimator $\hat{B}_i(t)$ is asymptotically design unbiased for $B_i(t)$ whereas $\hat{B}_{ic}(t)$ is asymptotically and conditionally design unbiased for $B_i(t)$ given

$i \in s$. The estimator $\hat{F}_{rkm}(t)$ is more complex than $\hat{F}^*(t)$ because it involves second-order inclusion probabilities otherwise it shares the same properties of the estimator $\hat{F}^*(t)$. Under model misspecification, $\hat{F}_{rkm}(t)$ performs better than \hat{F}_{cd} .

The asymptotic design variance of $\hat{F}_{rkm}(t)$ and an estimator of the variance was provided by Rao et al. (1990) as follows:

$$\begin{aligned} Var\{\hat{F}_{rkm}(t)\} &= \frac{1}{N^2} \frac{1}{2} \sum_{i \neq j} \sum_{j \in U} (\pi_i \pi_j - \pi_{ij}) \\ &\quad \left(\frac{I(y_i \leq t) - B_i(t)}{\pi_i} - \frac{I(y_j \leq t) - B_j(t)}{\pi_j} \right)^2 \end{aligned} \quad (23.2.27)$$

$$\begin{aligned} \hat{V}\{\hat{F}_{rkm}(t)\} &= \frac{1}{N^2} \frac{1}{2} \sum_{i \neq j} \sum_{j \in s} \frac{\pi_i \pi_j - \pi_{ij}}{\pi_{ij}} \\ &\quad \left(\frac{I(y_i \leq t) - \hat{B}_{ic}^*(j)}{\pi_i} - \frac{I(y_j \leq t) - \hat{B}_{jc}^*(i)}{\pi_j} \right)^2 \end{aligned} \quad (23.2.28)$$

where $\hat{B}_{ic}^*(j) = \sum_{k \in s} \frac{\pi_{ij}}{\pi_{ijk}} I\left(\frac{y_k - \hat{\beta}x_k}{\sqrt{v(x_k)}} \leq \frac{t - \hat{\beta}x_i}{\sqrt{v(x_i)}}\right) / \left(\sum_{k \in s} \frac{\pi_{ij}}{\pi_{ijk}}\right)$ and $\pi_{ijk}(i \neq j \neq k)$ is the inclusion probability for the i th, j th, and k th unit in the sample.

Dorfman (2009) proposed a modified estimator of $\hat{F}_{rkm}(t)$ as follows:

$$\begin{aligned} \hat{F}_{dm}(t) &= \frac{1}{N} \left[\sum_{i \in s} I(y_i \leq t) + \sum_{i \in U-s} \left\{ \frac{1}{n} \sum_{j \in s} I\left(\frac{y_i - \hat{\beta}x_j}{\sqrt{v(x_j)}} \leq \frac{t - \hat{\beta}x_i}{\sqrt{v(x_i)}}\right) \right\} \right. \\ &\quad \left. + \sum_{i \in s} \left(\frac{1}{\pi_i} - 1\right) R_i \right] \\ &= \hat{F}_{cd}(t) + \frac{1}{N} \sum_{i \in s} \left(\frac{1}{\pi_i} - 1\right) R_i \end{aligned} \quad (23.2.29)$$

where $R_i = I(y_i \leq t) - \frac{1}{n} \sum_{j \in s} I\left(\frac{y_j - \hat{\beta}x_j}{\sqrt{v(x_j)}} \leq \frac{t - \hat{\beta}x_i}{\sqrt{v(x_i)}}\right)$.

Both the estimator $\hat{F}_{dm}(t)$ and $\hat{F}_{rkm}(t)$ are equally efficient (Dorfman, 2009). $\hat{F}_{dm}(t)$ has additional advantage as it does not require computation of second-order inclusion probabilities. Wang and Dorfman (1996) proposed an alternative estimator $\hat{F}_{wdr}(t) = w\hat{F}_{dm}(t) + (1-w)\hat{F}_{rkm}(t)$, which is the weighted average of $\hat{F}_{dm}(t)$ and $\hat{F}_{rkm}(t)$. Simulation studies showed that the estimator $\hat{F}_{wdr}(t)$ performs better than $\hat{F}_{dm}(t)$ and $\hat{F}_{rkm}(t)$. Mak and Kuk (1993) proposed a modification of $\hat{F}_{cd}(t)$ estimator as follows:

$$\hat{F}_{mk}(t) = \frac{1}{N} \left[\sum_{i \in s} I(y_i \leq t) + \sum_{i \in U-s} \left\{ \frac{1}{n} \sum_{j \in s} \Phi \left(\frac{t - \hat{\beta}x_i}{\hat{\sigma} \sqrt{\nu(x_i)}} \right) \right\} \right]$$

where $\Phi \left(\frac{t - \hat{\beta}x_i}{\hat{\sigma} \sqrt{\nu(x_i)}} \right)$ is the standard normal distribution function and $\hat{\sigma}^2$ is the weighted list square estimate of σ^2 obtained from the model (23.2.17). The main advantage of using $\hat{F}_{mk}(t)$ over $\hat{F}_{cd}(t)$ is ease of computation but their relative performances are not known.

23.2.5 Nonparametric Regression Method

For the model-based and model-assisted estimators, we used parametric regression where a single smooth function $\mu(x)$ was estimated over the entire possible range of the auxiliary variable x . In a nonparametric regression we approximate $\mu(x)$ locally by placing more weight on y_i 's corresponding to x_i , which are close to x . In particular, suppose we wish to approximate $\mu(x)$ by a polynomial of degrees p viz.

$$\mu_p(x) = \beta_0 + \beta_1(x_i - x) + \cdots + \beta_p(x_i - x)^p \quad (23.2.30)$$

then the regression coefficients $\beta_0, \beta_1, \dots, \beta_p$ may be estimated by using weighted least square method using kernel weight

$$K \left(\frac{x_i - x}{b} \right)$$

attached to the i th observation. The kernel function $K(u)$ is a symmetric positive function of u , which decreases as $|u|$ increases. The parameter $b(>0)$ is the smoothing parameter, which is referred to bandwidth. Thus the normal equations are obtained from the selected sample by minimizing

$$\sum_{i \in s} (y_i - \beta_0 - \beta_1(x_i - x) - \cdots - \beta_p(x_i - x))^2 \left\{ \frac{1}{\pi_i} K \left(\frac{x_i - x}{b} \right) \right\}$$

with respect to $\beta_0, \beta_1, \dots, \beta_p$. Here we assume, without loss of generality, that the selected sample s contains first n units of the population. The estimates of $\boldsymbol{\beta}' = (\beta_0, \beta_1, \dots, \beta_p)$ are given by

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'_x \mathbf{W}_x \mathbf{X}_x)^{-1} \mathbf{X}'_x \mathbf{W}_x \mathbf{y}_s \quad (23.2.31)$$

where

$$\mathbf{X}_x = \begin{pmatrix} 1 & x_1 - x & \dots & (x_1 - x)^p \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ 1 & x_n - x & \cdot & (x_n - x)^p \end{pmatrix},$$

$$\mathbf{W}_x = \text{diag} \left(\frac{1}{\pi_1} K \left(\frac{x_1 - x}{b} \right) \dots \frac{1}{\pi_n} K \left(\frac{x_n - x}{b} \right) \right) \text{ and } \mathbf{y}'_s = (y_1, \dots, y_n).$$

The predicted value of y at the point $x = x_i$ is the intercept. So we write

$$\hat{y}_i = \hat{\mu}_p(x_i) = \hat{\beta}_0 = \mathbf{e}'_1 \left(\mathbf{X}'_{x_i} \mathbf{W}_{x_i} \mathbf{X}_{x_i} \right)^{-1} \mathbf{X}'_{x_i} \mathbf{W}_{x_i} \mathbf{y} \quad (23.2.32)$$

where \mathbf{e}_1 is a $(p+1)$ column vector with 1 as the first element and 0 elsewhere.

23.2.5.1 Nandaraya–Watson Estimator

For $p = 0$, Eq. (23.2.32) reduces to Nandaraya–Watson (1964) estimator:

$$\hat{\mu}_p(x_i) = \frac{\sum_{i \in s} K \left(\frac{x_i - x}{b} \right) \frac{y_i}{\pi_i}}{\sum_{i \in s} K \left(\frac{x_i - x}{b} \right) \frac{1}{\pi_i}} \quad (23.2.33)$$

The data analyst must choose p and b suitably. For further details, readers are referred to Ruppert et al. (2003).

23.2.5.2 Breidt and Opsomer Estimator

The local polynomial estimator for the population total is given by

$$\hat{Y}_{bo} = \sum_{i \in s} \frac{y_i - \hat{\mu}_p(x_i)}{\pi_i} + \sum_{i \in U} \hat{\mu}_p(x_i) \quad (23.2.34)$$

Now replacing y_i with $I(y_i \leq t)$ in Eq. (23.2.34), we get the Breidt and Opsomer (2000) estimator of the distribution function as

$$\hat{F}_p(t) = \frac{1}{N} \sum_{i \in s} \frac{I(y_i \leq t) - \hat{\mu}_p^*(x_i, t)}{\pi_i} + \frac{1}{N} \sum_{i \in U} \hat{\mu}_p^*(x_i, t) \quad (23.2.35)$$

where $\hat{\mu}_{lp}^*(x_i, t) = \mathbf{e}_1' \left(\mathbf{X}_{x_i}' \mathbf{W}_{x_i} \mathbf{X}_{x_i} \right)^{-1} \mathbf{X}_{x_i}' \mathbf{W}_{x_i} \mathbf{I}_{st}$ and $\mathbf{I}_{st} = (I(y_1 \leq t), \dots, I(y_i \leq t), \dots, I(y_n \leq t))'$.

From Breidt and Opsomer (2000), we see that $\hat{F}_{lp}(t)$ maintains design consistency and asymptotic design-unbiased properties. The estimated variance of $\hat{F}_{lp}(t)$ was given by Johnson (2003) as

$$\hat{V}(\hat{F}_{lp}(t)) = \frac{1}{2} \sum_{i \neq j} \frac{\pi_{ij} - \pi_i \pi_j}{\pi_{ij} \pi_i \pi_j} \left(I(y_i \leq t) - \hat{\mu}_{lp}^*(x_i, t) \right) \left(I(y_j \leq t) - \hat{\mu}_{lp}^*(x_j, t) \right) \quad (23.2.36)$$

23.2.5.3 Kuo Estimator

Kuo (1988) proposed the following estimator for the distribution function $F(t)$:

$$\hat{F}_{ko}(t) = \frac{1}{N} \left[\sum_{i \in s} I(y_i \leq t) + \sum_{j \in U-s} \sum_{i \in s} w_{ij} I(y_i \leq t) \right] \quad (23.2.37)$$

where $w_{ij} = \frac{K\left(\frac{x_j - x_i}{b}\right)}{\sum_{i \in s} K\left(\frac{x_j - x_i}{b}\right)}$, $K(z) = e^{-z^2/2}$ is the standard normal density (kernel).

Dorfman and Hall (1993) provided the expressions of asymptotic bias and variance of $\hat{F}_{ko}(t)$.

23.2.5.4 Kuk Estimator

Kuk (1993) nonparametric regression estimator for $F(t)$ is given by

$$\hat{F}_{kk}(t) = \frac{1}{N} \sum_{j \in U} \hat{R}_j(t) \quad (23.2.38)$$

where

$$\hat{R}_j(t) = \frac{\sum_{i \in s} \frac{1}{\pi_i} w\left(\frac{x_j - x_i}{b}\right) W\left(\frac{t - y_i}{b}\right)}{\sum_{i \in s} \frac{1}{\pi_i} w\left(\frac{x_j - x_i}{b}\right)}, \quad W(z) = e^z / (1 + e^z) \text{ is the}$$

standard logistic distribution function with density $w(z) = e^z / (1 + e^z)^2$, and b is the band with parameter used to control the amount of smoothing.

It should be noted that both the estimators $\hat{F}_{ko}(t)$ and $\hat{F}_{kk}(t)$ meet the properties of the distribution function given in Eq. (23.2.1). For more

details, readers are referred to Kuo (1988), Kuk (1993), Dorfman (2009), among others. Dorfman and Hall (1993) considered a design-adjusted version of Kuo (1988) and Rao–Kovar–Mantel (1990) estimators while Silva and Skinner (1995) suggested a poststratified estimator of $F(t)$.

23.2.6 Calibration Method

In Deville and Särndal's (1992) calibration method, the Horvitz–Thompson estimator $\hat{Y}_{ht} = \frac{1}{N} \sum_{i \in s} d_i y_i$ with $d_i = 1/\pi_i$ was calibrated as

$$\hat{Y}_c = \sum_{i \in s} w_i y_i$$

The weights w_i 's were chosen to minimize the distance

$$\Phi = \sum_{i \in s} \frac{(w_i - d_i)^2}{d_i q_i} \quad (23.2.39)$$

subject to the calibrating constraints $\sum_{i \in s} w_i x_i = X$. Here x_i 's are the values of the auxiliary variable x with known total X and q_i are suitably chosen weights. Details have been given in Section 9.7. Suppose that the study variable y is related to the auxiliary variable through the following superpopulation model

$$y_i = \mu(x_i, \theta) + \sqrt{v(x_i)} \epsilon_i$$

where θ is an unknown model parameter and ϵ_i 's are independently identically distributed with mean zero and variance σ^2 .

Let $\hat{\theta}$ be a suitable estimator of θ obtained from the selected sample using some standard procedure so that $\hat{\mu}_i = \mu(x_i, \hat{\theta})$ is an estimator for $\mu_i = \mu(x_i, \theta)$. Minimizing Eq. (23.2.39) subject to the calibrating constraints (i) $\sum_{i \in s} w_i = N$ and (ii) $\sum_{i \in s} w_i \hat{\mu}_i = \sum_{i \in U} \hat{\mu}_i$, Wu and Sitter (2001) derived the calibrated estimator for the population total Y as

$$\hat{Y}_c = \hat{Y}_{ht} - \hat{B} \left(\frac{1}{N} \sum_{i \in s} \frac{\hat{\mu}_i}{\pi_i} - \frac{1}{N} \sum_{i \in U} \hat{\mu}_i \right) \quad (23.2.40)$$

where $\hat{B} = \left(\sum_{i \in s} d_i q_i (\hat{\mu}_i - \bar{\mu}) (y_i - \bar{y}) \right) \left(\sum_{i \in s} d_i q_i (\hat{\mu}_i - \bar{\mu})^2 \right)^{-1}$,

$$\bar{y} = \sum_{i \in s} d_i q_i y_i / \sum_{i \in s} d_i q_i, \text{ and } \bar{\mu} = \sum_{i \in s} d_i q_i \hat{\mu}_i / \sum_{i \in s} d_i q_i.$$

For estimating distribution function $F(t)$, Wu and Sitter (2001) replaced y_i by $I(y_i \leq t)$, $\hat{\mu}_i$ by

$$\hat{G}_i = \frac{1}{n} \sum_{j \in s} I \left(\frac{y_j - \hat{\mu}_j}{\sqrt{v(x_j)}} \leq \frac{t - \hat{\mu}_i}{\sqrt{v(x_i)}} \right) \text{ and substituted } q_i = 1 \text{ in the}$$

expression (23.2.40). The resultant model calibrated estimator was $F(t)$ obtained as

$$\hat{F}_{us}(t) = \hat{F}_{ht}(t) - \hat{B}_{us} \left(\frac{1}{N} \sum_{i \in s} \frac{\hat{G}_i}{\pi_i} - \frac{1}{N} \sum_{i \in U} \hat{G}_i \right) \quad (23.2.41)$$

$$\text{where } \hat{F}_{ht}(t) = \frac{1}{N} \sum_{i \in s} \frac{I(y_i \leq t)}{\pi_i}, \quad \hat{B}_{us} = \frac{\left(\sum_{i \in s} d_i (\hat{G}_i - \bar{G}) (I(y_i \leq t) - \bar{I}) \right)}{\left(\sum_{i \in s} d_i (\hat{G}_i - \bar{G})^2 \right)},$$

$$\bar{G} = \frac{\sum_{i \in s} d_i \hat{G}_i}{\sum_{i \in s} d_i}, \text{ and } \bar{I} = \frac{\sum_{i \in s} d_i I(y_i \leq t)}{\sum_{i \in s} d_i}.$$

Wu and Sitter (2001) also proposed pseudoempirical likelihood estimator of $F(t)$ as

$$F_{us}^*(t) = \sum_{i \in s} \hat{p}_i I(y_i \leq t) \quad (23.2.42)$$

where the weights \hat{p}_i are obtained by maximizing pseudoempirical likelihood

$$l(p) = \sum_{i \in s} p_i I(y_i \leq t)$$

subject to (i) $p_i \geq 0$, (ii) $\sum_{i \in s} p_i = 1$, and (iii) $\sum_{i \in s} p_i \hat{G}_i = \frac{1}{N} \sum_{i \in U} \hat{G}_i$.

Since $p_i \geq 0$, $F_{us}^*(t)$ is a genuine distribution function. The calibrated estimator \hat{F}_{us} and the pseudocalibrated estimator F_{us}^* are asymptotically equivalent.

Rueda et al. (2007) considered a vector of auxiliary variables $\mathbf{x}_i = (x_{i1}, \dots, x_{ij}, \dots, x_{iq})'$, which are known for $i = 1, \dots, N$. They defined a pseudovariable $g_i = \hat{\beta} \mathbf{x}_i$ for $i = 1, \dots, N$, where

$\hat{\beta} = \left(\sum_{i \in s} d_i q_i \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \sum_{i \in s} d_i q_i \mathbf{x}_i y_i$ and q_i are known positive constants unrelated to $d_i = 1/\pi_i$. The proposed calibrated estimator is

$$\hat{F}_r = \frac{1}{N} \sum_{i \in s} w_i I(y_i \leq t) \quad (23.2.43)$$

where the weights w_i were derived by minimizing the chi-square distance

$$\Phi_s = \sum_{i \in s} \frac{(w_i - d_i)^2}{d_i q_i} \text{ subject to the constraints}$$

$$\frac{1}{N} \sum_{i \in s} w_i I(g_i \leq t_j) = F_g(t_j) = \frac{1}{N} \sum_{i \in U} I(g_i \leq t_j), j = 1, \dots, p$$

for suitably chosen p points $t_1 < t_2 < \dots < t_p$. Rueda et al. (2007), showed by simulation studies based on the actual data, that their estimator \hat{F}_r performs similar to \hat{F}_{cd} but better than \hat{F}_{rkm} .

23.2.7 Method of Poststratification

Silva and Skinner (1995) poststratified the initial sample s into H strata $U_1, \dots, U_h, \dots, U_H$ where the unit $i \in U_h$ if $x_{(h-1)} < x_i \leq x_{(h)}$ for $h = 1, \dots, H$, $-\infty = x_{(0)} < x_{(1)} < x_{(2)} < \dots < x_{(H)} = \infty$. Let s_h be the sample of size n_h from the h th stratum and $n = \sum_{h=1}^H n_h$. Silva and Skinner (1995) assumed that N_h , the stratum size of U_h , is known and n is so large that the probability of $n_h = 0$ is zero. Under this assumption, they proposed the following estimator of $F(t)$ as

$$\hat{F}_{ss}(t) = \sum_{h=1}^H W_h \hat{F}_h^*(t) \quad (23.2.44)$$

where $W_h = \frac{N_h}{N}$ and $\hat{F}_h^*(t) = \left(\sum_{i \in s_h} \frac{I(y_i \leq t)}{\pi_i} \right) / \left(\sum_{i \in s_h} \frac{1}{\pi_i} \right)$

It can be easily checked that the estimator $\hat{F}_{ss}(t)$ possesses all the properties of the distribution function stated in Eq. (23.2.1). Following Rao et al. (1990), Silva and Skinner (1995) derived the approximate expression of variance of $\hat{F}_{ss}(t)$ as

$$\begin{aligned} Var\{\hat{F}_{ss}(t)\} &\cong \frac{1}{N^2} \frac{1}{2} \sum_{i \neq j} \sum_{j \in U} (\pi_i \pi_j - \pi_{ij}) \\ &\times \left(\frac{I(y_i \leq t) - F_{h(i)}^*(t)}{\pi_i} - \frac{I(y_j \leq t) - F_{h(j)}^*(t)}{\pi_j} \right)^2 \end{aligned} \quad (23.2.45)$$

where $h(k)$ is the poststratum to which the unit k belongs and $F_h^*(t) = \frac{1}{N_h} \sum_{i \in U_h} I(y_i \leq t)$ is the population distribution function of the h th stratum $h = 1, \dots, H$.

An estimator of variance of $\hat{F}_{ss}(t)$ was presented by Silva and Skinner (1995) as

$$\hat{V}ar\{\hat{F}_{ss}(t)\} \cong \frac{1}{N^2} \frac{1}{2} \sum_{i \neq j} \sum_{j \in s} \frac{(\pi_i \pi_j - \pi_{ij})}{\pi_{ij}} \left(\frac{I(y_i \leq t) - \hat{F}_{h(i)}^*}{\pi_i} - \frac{I(y_j \leq t) - \hat{F}_{h(j)}^*(t)}{\pi_j} \right)^2 \quad (23.2.46)$$

where $\hat{F}_{h(i)}^*$ is an estimate of $F_{h(i)}^*$.

23.2.8 Empirical Comparison of the Estimators

Rao et al. (1990) did comprehensive studies of the performances of the design-based estimators $\hat{F}_h(t)$, $\hat{F}_R(t)$, $\hat{F}_D(t)$, $\hat{F}^*(t)$ and the model-based estimator $\hat{F}_{cd}(t)$ based on a population of sugarcane farms considered by Chambers and Dunstan (1986). The gross value of cane was treated as the study variable y while area under sugarcane was taken as auxiliary variable x . The population was found to obey the model (23.2.17) with $\nu(x) = \sqrt{x}$. Under simple random sampling, the relative biases of $\hat{F}_D(t)$ and $\hat{F}_{rkm}(t)$ are negligible and much less than that of $\hat{F}_R(t)$. Whence efficiency is concerned, $\hat{F}_{rkm}(t)$ is more efficient than $\hat{F}_D(t)$ and $\hat{F}(t)$. The ratio estimator $\hat{F}_R(t)$ is less efficient than $\hat{F}_D(t)$. The model-based estimator $\hat{F}_{cd}(t)$ was found to be more efficient than the design-based estimators possibly because the population obeys the model. Rao et al. (1990) also considered the Hansen et al. (1983) population with $N = 14,000$ units to study the effect of model misspecification and computed relative bias and relative root mean square errors of the estimators $\hat{F}_h(t)$, $\hat{F}_R(t)$, $\hat{F}_D(t)$, $\hat{F}_{cd}(t)$, and $\hat{F}_{rkm}(t)$. It was found that design-based estimators have negligible bias, even less than 1%, whereas the model-based estimator $\hat{F}_{cd}(t)$ has a much larger relative bias of about 20%. The design-based estimator $\hat{F}_{rkm}(t)$ was found to be much more efficient than the other design-based estimators $\hat{F}_R(t)$ and $\hat{F}_D(t)$ whereas the model-based estimator $\hat{F}_{cd}(t)$ was found least efficient.

Silva and Skinner (1995) conducted simulation studies to compare performances of the estimators $\hat{F}_h(t)$, $\hat{F}_{ss}(t)$, $\hat{F}_{ko}(t)$, $\hat{F}_{kk}(t)$, $\hat{F}_{cd}(t)$, and $\hat{F}_{rlm}(t)$ based on two populations. The first population was the sugarcane farm data mentioned above, which was originally considered by Chambers and Dunstan (1986) and later Rao et al. (1990) and Kuk (1993). The second population comprises of 430 farms with 50 or more beef cattle, which was originally used by Chambers et al. (1993) and then Kuk (1993). Population 1 is a good working model for Eq. (23.2.17) as stated earlier whereas Population 2 is not. As far as relative bias is concerned, for both the populations, $\hat{F}_{ss}(t)$ was found to have least bias while $\hat{F}_{rlm}(t)$ occupied second place. $\hat{F}_{ko}(t)$, $\hat{F}_{kk}(t)$, and $\hat{F}_{cd}(t)$ were found to have relatively high absolute relative biases. For Population 1, $\hat{F}_{cd}(t)$ had least mean squared error (MSE) as expected because the model fits the population. The second and third places were occupied by $\hat{F}_{kk}(t)$ and $\hat{F}_{rlm}(t)$, respectively. Surprisingly, for Population 2, $\hat{F}_{cd}(t)$ also had the smallest MSE followed by $\hat{F}_{kk}(t)$. The estimators $\hat{F}_{ss}(t)$ and $\hat{F}_{ko}(t)$ were found to have moderate MSE. For further details, interested readers are referred to Rao et al. (1990) and Silva and Skinner (1995).

23.3 ESTIMATION OF QUANTILES

The α th quantile $\theta_y(\alpha)$, $0 < \alpha < 1$ of a finite population vector $\mathbf{y} = (y_1, \dots, y_N)$ is defined as

$$\theta_y(\alpha) = \inf\{t : F_y(t) \geq \alpha\} \quad (23.3.1)$$

where $F_y(t)$ is the distribution function of y . In case $\hat{F}_y(t)$, an estimator of $F_y(t)$, is a monotonic nondecreasing function of t , the customary estimator of $\theta_y(\alpha)$ is obtained as

$$\hat{\theta}_y(\alpha) = \inf\{t : \hat{F}_y(t) \geq \alpha\} \quad (23.3.2)$$

Let $\hat{F}_x(t)$ be the customary estimator of $F_x(t)$. In case the population α th quantile $\theta_x(\alpha)$ of x is known, the ratio estimator of $\theta_y(\alpha)$ is given by

$$\hat{\theta}_{ry}(\alpha) = \frac{\hat{\theta}_y(\alpha)}{\hat{\theta}_x(\alpha)} \theta_x(\alpha) \quad (23.3.3)$$

Similarly, a difference estimator of $\theta_y(\alpha)$ is given by

$$\hat{\theta}_{dy}(\alpha) = \hat{\theta}_y(\alpha) - \hat{R} \left\{ \hat{\theta}_x(\alpha) - \theta_x(\alpha) \right\} \quad (23.3.4)$$

where $\hat{R} = \frac{\sum_{i \in s} y_i / \pi_i}{\sum_{i \in s} x_i / \pi_i}$ is a consistent estimator of the population ratio

$$R = Y/X.$$

Both the estimators $\hat{\theta}_{ry}(\alpha)$ and $\hat{\theta}_{dy}(\alpha)$ reduce to $\theta_y(\alpha)$ if $y_i \propto x_i \forall i \in U$. In this case the variances of the estimators become zero. Hence the estimators $\hat{\theta}_{ry}(\alpha)$ and $\hat{\theta}_{dy}(\alpha)$ are expected to produce a considerable gain in efficiency over $\hat{\theta}_y(\alpha)$ if y_i is approximately proportional to x_i . Rao et al. (1990) derived the variances of $\hat{\theta}_{ry}(\alpha)$ and $\hat{\theta}_{dy}(\alpha)$ and also their unbiased estimators. These are omitted here because of their complexities.

23.4 ESTIMATION OF MEDIAN

The median of a variable y is obtained by substituting $\alpha = 1/2$ in Eq. (23.3.1). Thus the population median of y and x are, respectively,

$$\begin{aligned}\tilde{\mu}_y &= \theta_y(1/2) = \inf\{t : F_y(t) \geq 1/2\} \text{ and} \\ \tilde{\mu}_x &= \theta_x(1/2) = \inf\{t : F_x(t) \geq 1/2\}\end{aligned}\quad (23.4.1)$$

If the population median $\tilde{\mu}_x = \theta_x(1/2)$ of the auxiliary variable x is known, then the ratio estimator of the population median $\tilde{\mu}_y$ is obtained from Eq. (23.3.3) as

$$\hat{\tilde{\mu}}_{ry} = \hat{\theta}_{ry}(1/2) = \frac{\hat{\tilde{\mu}}_y}{\hat{\tilde{\mu}}_x} \tilde{\mu}_x \quad (23.4.2)$$

where $\hat{\tilde{\mu}}_x$ and $\hat{\tilde{\mu}}_y$ are the sample medians of x and y , respectively.

Kuk and Mak (1989) proposed the following modifications of the ratio estimator of $\hat{\tilde{\mu}}_{ry}$ based on simple random sampling.

23.4.1 Position Estimator and Stratification Estimator

At first, let us arrange y_i 's, $i \in s$ in order of magnitude as $y_{(1)} \leq y_{(2)} \leq \dots \leq y_{(n)}$. Let i_0 be the number of observations of y_i 's less than or equal to the population median $\tilde{\mu}_y$, i.e., $y_{(i_0)} \leq \tilde{\mu}_y < y_{(i_0+1)}$. Clearly, i_0 is unknown because $\tilde{\mu}_y$ is unknown. Let $p = i_0/n$, then $\tilde{\mu}_y$ is approximately the sample p th quantile $\hat{\theta}_y(p)$. Suppose \hat{p} is an estimator of p , then $\hat{\theta}_y(\hat{p})$ is an estimator of the median $\tilde{\mu}_y$. $\hat{\tilde{\mu}}_y$, the sample median of y can be viewed as the special estimator $\hat{\theta}_y(\hat{p})$ with $\hat{p} = 1/2$. To estimate \hat{p} , Kuk and Mak (1989) considered the following two-way classified table.

	$x \leq \tilde{\mu}_x$	$x > \tilde{\mu}_x$	Total
$y \leq \tilde{\mu}_y$	P_{11}	P_{12}	P_{10}
$y > \tilde{\mu}_y$	P_{21}	P_{22}	P_{20}
Total	P_{01}	P_{02}	1

In the table above, P_{ij} denotes the proportion of units in the population that belongs to the (i, j) th cell. Let n_x be the number of observations in the sample s with x_i 's less than or equal to $\tilde{\mu}_x$. In case P_{ij} are known, an estimate of p is given by

$$\begin{aligned}\hat{p} &= \frac{1}{n} \left(n_x \frac{P_{11}}{P_{01}} + (n - n_x) \frac{P_{12}}{P_{02}} \right) \\ &\cong \frac{2}{n} \left(n_x P_{11} + (n - n_x) \left(\frac{1}{2} - P_{11} \right) \right)\end{aligned}\quad (23.4.3)$$

where $P_{0j} = P_{1j} + P_{2j} \cong 1/2$ for $j = 1, 2$.

In practice, P_{ij} are generally unknown. So, we estimate P_{ij} by p_{ij} , the proportion of units in the sample that fall in the (i, j) th cell, i.e., p_{11} is the proportion of observation in the class $x \leq \hat{\mu}_x$ (=sample median of x) and $y \leq \hat{\mu}_y$. Now replacing P_{11} by p_{11} in Eq. (23.4.3), we find an estimate of p as

$$\hat{p}_1 \cong \frac{2}{n} \left(n_x p_{11} + (n - n_x) \left(\frac{1}{2} - p_{11} \right) \right) \quad (23.4.4)$$

Thus an estimator of $\tilde{\mu}_y$ is given by

$$\hat{\mu}_{py} = \hat{\theta}_y(\hat{p}_1) \quad (23.4.5)$$

The estimator $\hat{\mu}_{py}$ was termed as the “position estimator” by Kuk and Mak (1989).

For a given value of $y = t$, let $\tilde{F}_{1y}(t)$ be the proportion of those units in the sample with $x \leq \tilde{\mu}_x$, which have y -values less than equal to t and $\tilde{F}_{2y}(t)$ be the proportion of those units in the sample with $x > \tilde{\mu}_x$, which have y -values less than equal to t . Then $F_y(t)$ can be estimated as

$$\begin{aligned}\tilde{F}_y(t) &= \frac{1}{N} [N_x \tilde{F}_{1y}(t) + (N - N_x) \tilde{F}_{2y}(t)] \\ &\cong \frac{1}{2} [\tilde{F}_{1y}(t) + \tilde{F}_{2y}(t)]\end{aligned}\quad (23.4.6)$$

where N_x is the number of units in the population with $x \leq \tilde{\mu}_x$.

Since $\widetilde{F}_y(t)$ is a distribution function, an estimator of the median $\widetilde{\mu}_y$ is obtained as

$$\widehat{\mu}_{sy} = \inf \left\{ \gamma : \widetilde{F}_y(t) \geq \frac{1}{2} \right\} \quad (23.4.7)$$

Kuk and Mak (1989) called the estimator $\widehat{\mu}_{sy}$ as the “stratification estimator.”

23.4.2 Comparison of the Efficiencies

Gross (1980) and Kuk and Mak (1989) derived the asymptotic distributions of the median estimators $\widehat{\mu}_y$, $\widehat{\mu}_{ry}$, $\widehat{\mu}_{py}$, and $\widehat{\mu}_{sy}$ when $N \rightarrow \infty$, $n \rightarrow \infty$ and $n/N \rightarrow f$, $0 \leq f \leq 1$. They assumed that as $N \rightarrow \infty$, the bivariate distribution of (x, y) approaches a continuous distribution with marginal densities $f_x(x)$ and $f_y(y)$ for x and y , respectively. Gross (1980) proved that the sample median $\widehat{\mu}_y$ is consistent and asymptotically normally distributed with mean $\widetilde{\mu}_y$ and variance

$$Var(\widehat{\mu}_y) = \frac{(1-f)}{4n} \left\{ f_y(\widetilde{\mu}_y) \right\}^{-2} \quad (23.4.8)$$

Kuk and Mak (1989) derived the following results:

(i) $\widehat{\mu}_{ry}$ is asymptotically normal with mean $\widetilde{\mu}_y$ and variance

$$Var(\widehat{\mu}_{ry}) = \frac{(1-f)}{n} \left[\frac{1}{4} \left(f_y(\widetilde{\mu}_y) \right)^{-2} + \frac{1}{4} \widetilde{R}^2 (f_x(\widetilde{\mu}_x))^{-2} - 2\widetilde{R} \left(\frac{f_y(\widetilde{\mu}_y)}{f_x(\widetilde{\mu}_x)} \right)^{-1} \left(P_{11} - \frac{1}{4} \right) \right] \quad (23.4.9)$$

where $\widetilde{R} = \widetilde{\mu}_y / \widetilde{\mu}_x$.

(ii) $\widehat{\mu}_{py}$ and $\widehat{\mu}_{sy}$ both asymptotically follow the same distribution that is normal with mean $\widetilde{\mu}_y$ and variance

$$Var(\widehat{\mu}_{py}) = Var(\widehat{\mu}_{sy}) = \frac{2(1-f)P_{11}(1-2P_{11})}{n} \left(f_y(\widetilde{\mu}_y) \right)^{-2} \quad (23.4.10)$$

where P_{11} is the proportion of units in the population with $x \leq \widetilde{\mu}_x$ and $y \leq \widetilde{\mu}_y$. The probability P_{11} can be regarded as a measure of concordance.

It is important to note that if $P_{11} = 1/2$, then the asymptotic variances $Var(\widehat{\mu}_{py})$ and $Var(\widehat{\mu}_{sy})$ are both equal to zero. The expressions (23.4.8) and (23.4.10) show that the estimators $\widehat{\mu}_{py}$ and $\widehat{\mu}_{sy}$ are asymptotically more efficient than the sample median because

$$Var(\widehat{\mu}_{sy}) - Var(\widehat{\mu}_y) = \frac{1}{4}(4P_{11} - 1)^2 \geq 0 \quad (23.4.11)$$

The estimator $\widehat{\mu}_{py}$ becomes asymptotically more efficient than the sample median $\widehat{\mu}_y$ if

$$\rho_c > \frac{1}{2} \left(\frac{\widetilde{R} f_y(\widetilde{\mu}_y)}{f_x(\widetilde{\mu}_x)} \right) \quad (23.4.12)$$

where $\rho_c = 4(P_{11} - 1/4)$ varies from -1 to 1 as P_{11} increases from 0 to $1/2$.

23.4.3 Further Generalization

Let the vector of the auxiliary variable $\mathbf{x} = (x_1, \dots, x_N)$ be known and the range of the auxiliary variable be partitioned into r mutually exclusive and exhaustive class intervals $(a_0, a_1], (a_1, a_2], \dots, (a_{r-1}, a_r]$ with $a_0 = 0$, $a_r = \infty$, and $a_j = \theta_x(\alpha_j)$. Let P_{1j} be the proportion of the units in the population with $y \leq \widetilde{\mu}_y$ and x falling in the class $(a_{j-1}, a_j]$ for $j = 1, \dots, r$ and let p_{1j} be the proportion of units in s with $y \leq \widehat{\mu}_y$ and x in $(\widehat{a}_{j-1}, \widehat{a}_j]$ where $\widehat{a}_j = \widehat{\theta}_x(\alpha_j)$. Also, let P_{0j} be the proportion of units in the population with $x \in (a_{j-1}, a_j]$ and p_{0j} be the proportion in the sample with $x \in (\widehat{a}_{j-1}, \widehat{a}_j]$, whereas n_{xj} is the number of units in the sample with $x \in (a_{j-1}, a_j]$.

Proportions in population

	$(a_0, a_1]$	—	$(a_{j-1}, a_j]$	—	$(a_{r-1}, a_r]$
$y \leq \widetilde{\mu}_y$	P_{11}	—	P_{1j}	—	P_{1r}
$y > \widetilde{\mu}_y$	P_{21}	—	P_{2j}	—	P_{2r}
Total	P_{01}	—	P_{0j}	—	P_{0r}

Proportions in sample

	$(a_0, a_1]$	—	$(a_{j-1}, a_j]$	—	$(a_{r-1}, a_r]$
$y \leq \widehat{\mu}_y$	p_{11}	—	p_{1j}	—	p_{1r}
$y > \widehat{\mu}_y$	p_{21}	—	p_{2j}	—	p_{2r}
Total	p_{01}	—	p_{0j}	—	p_{0r}

The estimated proportion of γ 's in the sample that is less than equal to $\tilde{\mu}_\gamma$ is given by Kuk and Mak (1989) as

$$\begin{aligned}\tilde{p}_{g1} &= \frac{1}{n} \sum_{j=1}^r n_{xj} \frac{p_{1j}}{p_{0j}} \\ &= \frac{1}{n} \sum_{j=1}^r n_{xj} \frac{p_{1j}}{(\alpha_j - \alpha_{j-1})}\end{aligned}\quad (23.4.13)$$

Finally, the estimated populated median is

$$\hat{\mu}_{gpy} = \hat{\theta}_\gamma(\tilde{p}_{g1}) \quad (23.4.14)$$

Let $\hat{F}_{y\gamma}(t)$ be the proportion of the units in the sample with $x \in (a_{j-1}, a_j]$ that have γ -values less than or equal to t . Then $F_\gamma(t)$ may be estimated by

$$\hat{F}_\gamma(t) = \sum_{j=1}^r (\alpha_j - \alpha_{j-1}) \hat{F}_{y\gamma}(t)$$

Consequently the estimated median is

$$\hat{\mu}_{gsy} = \inf\{\gamma : \hat{F}_\gamma(t) \geq 1/2\} \quad (23.4.15)$$

Kuk and Mak (1989) showed that the asymptotic variances of $\hat{\mu}_{gpy}$ and $\hat{\mu}_{gsy}$ are the same and equal to

$$V_{ar}(\hat{\mu}_{gpy}) = V_{ar}(\hat{\mu}_{gsy}) = \frac{1-f}{n} \frac{1}{(f_\gamma(\tilde{\mu}_\gamma))^2} \left(\frac{1}{2} - \sum_{j=1}^r \frac{P_{1j}^2}{(\alpha_j - \alpha_{j-1})} \right) \quad (23.4.16)$$

23.4.4 Empirical Comparison

Kuk and Mak (1989) compared efficiencies of the estimators $\hat{\mu}_\gamma$, $\hat{\mu}_{\gamma\gamma}$, $\hat{\mu}_{py}$, and $\hat{\mu}_{sy}$ empirically using four populations named Hospitals, Counties 70, Villages, and Factories. The first two populations were used by Royall and Cumberland (1981), where x and γ are well correlated and a linearity relationship holds. For the population Villages (Murthy, 1967), where x (area in 1951) and γ (number of households in 1961) are poorly correlated, the probability of concordance P_{11} is also low. For the population Factories (Murthy, 1967), x (number of workers) and γ (output) are not linearly

related. From the first two populations, 500 independent samples of suitable sizes were selected by SRSWOR method while from each of the remaining other two populations, 1000 independent samples of suitable sizes were selected. Empirical studies reveal that for the populations Hospitals and Counties where x and y are linearly related, all three estimators $\hat{\mu}_{ry}$, $\hat{\mu}_{py}$, and $\hat{\mu}_{sy}$ have much lower mean square errors than that of the sample median $\hat{\mu}_y$. For the population Villages where x and y are poorly related, the ratio estimator $\hat{\mu}_{ry}$ performed worse than the sample median $\hat{\mu}_y$. But the other two estimators $\hat{\mu}_{py}$ and $\hat{\mu}_{sy}$ fare better than $\hat{\mu}_y$. For the nonlinear population Factories, $\hat{\mu}_{py}$ and $\hat{\mu}_{sy}$ still performed better than $\hat{\mu}_y$ where the ratio estimator $\hat{\mu}_{ry}$ is again outperformed by the sample median. Thus the efficiencies of the position estimator $\hat{\mu}_{py}$ and stratification estimator $\hat{\mu}_{sy}$ do not depend on the validity of linearity assumption and hence it is much safer to use, than the ratio estimator $\hat{\mu}_{ry}$.

23.5 CONFIDENCE INTERVAL FOR DISTRIBUTION FUNCTION AND QUANTILES

The conventional $(1 - \alpha)100\%$ confidence interval of the distribution function $F(t)$ is

$$\left(\hat{F}(t) - z_{\alpha/2} \sqrt{\hat{V}[\hat{F}(t)]}, \hat{F}(t) + z_{\alpha/2} \sqrt{\hat{V}[\hat{F}(t)]} \right) \quad (23.5.1)$$

where $z_{1-\alpha/2}$ is the $1 - \alpha/2$ quantile from $N(0, 1)$ and $\hat{V}[\hat{F}(t)]$ is an estimator of $V[\hat{F}(t)]$. The validity of the confidence interval obviously depends on the asymptotic normality of the distribution

$$\frac{\hat{F}(t) - F(t)}{\sqrt{\hat{V}[\hat{F}(t)]}}$$

which is justified when the sample size is large. However, for small to moderate sample size, because of range constraint $0 \leq \hat{F}(t) \leq 1$, the distribution of $\hat{F}(t)$ for t at large or small quantiles is usually not symmetric. For small and moderate sample size, the performance of the interval (23.5.1) is very often unsatisfactory; coverage probability is generally lower than the nominal value and two tail probabilities are unbalanced (Wu, 1999).

Chen and Wu (2002) proposed an alternative method. In this method a smooth and monotone function g is chosen so that the distribution of

$\widehat{W} = \widehat{W}(t) = g[\widehat{F}(t)]$ is better approximated by the normal distributions. Two such popular transformations are logit transformation and complementary log-log transformation: $\widehat{W} = \log\left(\frac{\widehat{F}(t)}{1 - \widehat{F}(t)}\right)$ and $\log(-\log\{\widehat{F}(t)\})$. Noting that \widehat{W} asymptotical normal with mean $W = W(t) = g[F(t)]$ and variance $V(\widehat{W}) = (g'\{F(t)\})^2 V\{\widehat{F}(t)\}$, we can find a $(1 - \alpha)100\%$ confidence interval of W as

$$\left(\widehat{W} - z_{\alpha/2} \sqrt{\widehat{V}(\widehat{W})}, \widehat{W} + z_{\alpha/2} \sqrt{\widehat{V}(\widehat{W})} \right)$$

where $\widehat{V}(\widehat{W}) = (g'\{\widehat{F}(t)\})^2 \widehat{V}\{\widehat{F}(t)\}$ is a suitable estimator of $V(\widehat{W})$.

Finally, transformed confidence interval for $F(t)$ is obtained as

$$\left(g^{-1} \left\{ \widehat{W} - z_{\alpha/2} \sqrt{\widehat{V}(\widehat{W})} \right\}, g^{-1} \left\{ \widehat{W} + z_{\alpha/2} \sqrt{\widehat{V}(\widehat{W})} \right\} \right) \quad (23.5.2)$$

On the basis of simulation studies, Wu (1999) showed that performance of the intervals (23.5.1) and (23.5.2) is similar when t is in the middle range of quantiles but Eq. (23.5.2) is much better when t is small or large quantiles.

Let $\widehat{\theta}(\alpha)$ be an estimator of the $\alpha(0 < \alpha < 1)$ th quantile $\theta(\alpha)$ obtained from the inversion of estimated distribution function $\widehat{F}(t)$, i.e., $\widehat{F}\{\widehat{\theta}(\alpha)\} = \alpha$. Furthermore, let $L(t)$ and $U(t)$ be the lower and upper $(1 - \alpha)100\%$ confidence intervals of $F(t)$ obtained from any of the formulas (23.5.1) and (23.5.2). Then following Woodruff (1952), the confidence interval of the $\alpha(0 < \alpha < 1)$ th quantile $\theta(\alpha)$ is obtained as

$$\left(\widehat{F}^{-1} \left[L\{\widehat{\theta}(\alpha)\} \right], \widehat{F}^{-1} \left[U\{\widehat{\theta}(\alpha)\} \right] \right) \quad (23.5.3)$$

Obviously the confidence interval based on the formula (23.5.2) is expected to perform better than that on Eq. (23.5.1).

Godambe and Thompson (1999) used estimating functions for determination of confidence intervals of distribution functions and quantiles while Chen and Sitter (1999), Wu and Sitter (2001), among others, used empirical likelihood methods. Details have been given in Chapters 22 and 25, respectively.

23.6 CONCLUDING REMARKS

In surveys, estimation of distribution function and quantiles often play an important role. Rao et al. (1990) showed how one can use auxiliary information x to construct customary design-based ratio $\hat{F}_R(t)$ and difference estimator $\hat{F}_D(t)$ of the population distribution function $F(t)$, which can improve upon $\hat{F}(t)$, the conventional estimator of the distribution function $F(t)$. The proposed estimators $\hat{F}_R(t)$ and $\hat{F}_D(t)$ are asymptotically design unbiased but not model unbiased. Chambers and Dunstan (1986) proposed a model-based estimator $\hat{F}_{cd}(t)$, which is efficient if the assumed model is valid, but this estimator can perform poorly under model misspecification. Rao et al. (1990) also proposed model-based estimators $\hat{F}^*(t)$ and $\hat{F}_{rkm}(t)$, which are both asymptotically design unbiased and model-unbiased under the model (23.2.17) with $\nu(x_i) = \sqrt{x_i}$. Moreover, $\hat{F}_{rkm}(t)$ can be modified using multiauxiliary variables. Poststratification estimator $\hat{F}_{ss}(t)$ was proposed by Silva and Skinner (1995) and it retains the property of asymptotic design unbiasedness. To compute Kuk's (1993) estimator $\hat{F}_{kk}(t)$, the values of the study variable y need to be scaled since the bandwidth that is used to control the smoothing is the same for both the variables x and y . The estimators $\hat{F}_R(t)$, $\hat{F}_D(t)$, $\hat{F}^*(t)$, $\hat{F}_{cd}(t)$, and $\hat{F}_{ss}(t)$ are advantageous because the variance estimators are readily available and easy to compute. Although the estimator of variance of $\hat{F}_{rkm}(t)$ is readily available, it can be extremely complex for a varying probability sampling scheme because of the involvement of third-order inclusion probabilities π_{ijk} that are difficult to compute. Extensive simulation studies were done by Chambers and Dunstan (1986), Rao et al. (1990), Chambers et al. (1993), Kuk (1993), and Silva and Skinner (1995), among others, to compare performances among the proposed estimators. The simulation studies do not finger out any particular estimator as the best in all situations. However, $\hat{F}_{rkm}(t)$ seems to perform well in most of the situations irrespective of the failure of the model and so it is safe to be used.

Estimators of the distribution functions can be extended to estimation of population quantiles and medians. Kuk and Mak (1989) proposed the alternative position estimator $\hat{\mu}_{py}$ and stratification estimator $\hat{\mu}_{sy}$ to estimate the median under SRSWOR sampling. Empirical studies of Kuk and Mak (1989) reveal that the estimators $\hat{\mu}_{py}$ and $\hat{\mu}_{sy}$ are more efficient than the sample median $\hat{\mu}_y$ and the ratio estimator $\hat{\mu}_{ry}$ even if the relation between x and y is not linear. The use of empirical likelihood method for estimating

finite population distribution function was considered by Owen (1988), Chen and Sitter (1999), Wu and Sitter (2001). The proposed estimators use auxiliary information effectively at the estimation stage and possess attractive properties. Godambe and Thompson (1999) used the method of estimating function whereas Chen and Wu (2002) used calibrated pseudoempirical likelihood methods for determining confidence intervals. Further discussions are given by Johnson (2003), Drofman (2009), among others.

23.7 EXERCISES

23.7.1 Let a sample of size 10 be selected from a population of size 50 by simple random sampling with replacement method.

Selected sample	1	2	3	4	5	6	7	8	9	10
y -values	10	8	6	15	3	8	5	6	3	1

- (i) Sketch the empirical distribution function.
- (ii) Estimate the 25th, 50th, and 80th percentiles from the graph. Estimate the standard errors of the estimators used. Determine 95% confidence interval of the population median.

23.7.2 A sample of size of size 8 is selected from a finite population of size 50 using Poisson sampling scheme. The following table gives the units selected in the sample, y -values, and inclusion probabilities (π_i) of the selected units.

Selected sample	1	2	3	4	5	6	7	8
y -values	15	10	5	4	6	15	20	15
π_i	0.10	0.15	0.20	0.25	0.15	0.10	0.25	0.10

Sketch the distribution function using the formulae.

$$(i) \hat{F}(t) = \frac{1}{N} \sum_{i \in s} \frac{I(y_i \leq t)}{\pi_i}$$

$$(ii) \hat{F}_h(t) = \sum_{i \in s} \frac{I(y_i \leq t)}{\pi_i} \bigg/ \sum_{i \in s} \frac{1}{\pi_i}$$

$$(iii) \hat{F}_c(t) = 1 - \frac{1}{N} \sum_{i \in s} \frac{1 - I(y_i \leq t)}{\pi_i}$$

Discuss the appropriateness of the formulae. From the sketches, obtain estimates of the population median. Compute standard errors of the estimators used.

23.7.3 A sample of 10 plants is selected at random from 40 plants of a garden. The following table gives the height (y) of 10 plants along with the diameter (x) of all the 40 plants in that garden. Fit distribution functions over the data using the following methods:

(i) Ratio method: $\widehat{F}_R(t) = \frac{1}{N} \widehat{\lambda} \sum_{i \in U} \frac{I(\widehat{y}_i \leq t)}{\pi_i}$

where $\widehat{\lambda} = \left(\frac{\sum_{i \in s} \frac{I(y_i \leq t)}{\pi_i}}{\sum_{i \in s} \frac{I(\widehat{y}_i \leq t)}{\pi_i}} \right)$ and $\widehat{y}_i = \left(\frac{\sum_{i \in s} y_i}{\sum_{i \in s} 1} \right) x_i$

(ii) Difference method:

$$\widehat{F}_D(t) = \frac{1}{N} \left[\sum_{i \in s} \frac{I(y_i \leq t)}{\pi_i} - \left\{ \sum_{i \in s} \frac{I(\widehat{y}_i \leq t)}{\pi_i} - \sum_{i \in U} I(\widehat{y}_i \leq t) \right\} \right]$$

From the fitted distribution functions, estimate median and 80th percentiles of the heights of the plants. Also estimate standard errors of the estimators used.

Plants	1	2	3	4	5	6	7	8	9	10
Diameter (in cm)	15	20	30	25	30	40	15	12	20	30
Height (in cm)	80	90	120	100	100					

Plants	11	12	13	14	15	16	17	18	19	20
Diameter (in cm)	15	20	30	25	30	40	15	12	20	30

Plants	21	22	23	24	25	26	27	28	29	30
Diameter (in cm)	10	20	18	28	50	60	45	20	25	30

Plants	31	32	33	34	35	36	37	38	39	40
Diameter (in cm)	15	25	20	10	10	30	25	40	60	30

23.7.4 Continuation of Exercise 27.7.3. Assume that y and x are related to the model $y_i = \beta x_i + \epsilon_i$, where ϵ_i 's are independent with $E_m(\epsilon_i) = 0$ and $V_m(\epsilon_i) = \sigma^2 x_i$. Estimate the median height of the plant using the Chamber and Dunstan estimator

$$\hat{F}_{cd}(t) = \frac{1}{N} \left[\sum_{i \in s} I(y_i \leq t) - \sum_{i \in U-s} \left\{ \frac{1}{n} \sum_{j \in s} I \left(\frac{(y_i - \hat{\beta} x_j)}{\sqrt{\nu(x_j)}} \leq \frac{(t - \hat{\beta} x_j)}{\sqrt{\nu(x_j)}} \right) \right\} \right]$$

where $\hat{\beta} = \sum_{j \in s} \frac{y_j x_j}{\sqrt{\nu(x_j)}} / \sum_{j \in s} \frac{x_j^2}{\sqrt{\nu(x_j)}}$ and $\nu(x_j) = \sigma^2 x_j^2$.

23.7.5 The following table relates to the daily wages in (US\$) of 40 factory workers selected at random from 120 workers.

120	80	160	200	400	550	300	400	500	150
140	90	150	180	300	250	200	180	150	250
270	380	920	150	200	225	420	100	120	150
80	75	150	200	500	590	250	140	150	200

Estimate the median wage of the factory workers and obtain 90% confidence interval of the median wages.