Sample Survey

- ② → Cluster sampling
- ③ → Two-stage Sampling
- ④ → Double Sampling
- ⑥ → PPS Sampling
- ⑤ → Sampling Design & Inference
- ① → Randomised Response

## Randomised Response Technique :

In some sample survey studies, some of the study variables are sensitive in the nature but they are legal to be asked. Examples of such study variables are

(i) addition to alcohol or drugs.

(ii) having a history of abortions at some area.

(iii) habit of gambling.

(iv) habit of being without ticket passenger etc.

Clearly, considering complete and correct information on such sensitive issues is very difficult for the interviewer. It is not a delied affair to bluntly ask a question (a stranger ~~question~~ question) about such sensitive and highly personal matter. Even if the interviewer doesn't refuse to answer such sensitive queries, we might wonder whether the replies provide by an interviewee are honest. This means that the collected data with the questionaires are very incomplete. In order to avoide excessive refusale or misleading responses to sensitive questions, we need new statistical methodology capable of (i) allowing the interviwee a satisfactory level of privacy and ~~the two~~ providing a basis which encourages greater coordination and

(ii) providing valid estimates of the population parameters under study.

Various statistical methods have been devised by the Statisticians dealing with such situations. All such techniques are known as "Randomized Response Techniques" the first study of these kinds were undertaken by Warner (1965) who developed a clever interviewing process design to reduce or eliminate non-sampling biases. in sample surveys at human population.

▣ Randomized Response Technique for one dichotomous qualitative sensitive characteristic :

---

>> What is the problem ?

Let Y be a qualitative characteristic with two versions A (sensitive) and $A^c$ (non-sensitive). For convenience, let

$$Y_i = \begin{cases} 0 & \text{, if the } i^{th} \text{ population unit have the non-sensitive version } A^c \\ 1 & \text{, otherwise.} \end{cases}$$

With respect to the characteristic, the population $U$ can be divided into two subpopulations $U_0$ & $U_1$ having the size $N_0$ and $N_1$ respectively. All individual in $U_1$ have the sensitive version of the characteristic and the goal of survey is to make a statistical study of the unknown population parameter $\phi = \dfrac{N_1}{N}$, which is the proportion of units in the population having the sensitive version of the characteristic Y.

Suppose we pose the following question to the $i^{th}$ population unit :

Do you belong to $U_1$ ?

Let, $Y_i$ be the corresponding response variable. Then it is not unreasonable to postulate the following modes of response

$$P(Y_i = 1 \mid i \in u_0) = 0$$

$$P(Y_i = 0 \mid i \in u_0) = 1$$

and

$$P(Y_i = 1 \mid i \in u_1) = q_{i1}$$

$$P(Y_i = 0 \mid i \in u_1) = q_{i2}$$

$$P(\text{No response from individual } i \mid i \in u_1)$$
$$= q_{i3}$$

In other words, if the $i$th unit belongs to $u_0$ (ie. the non-sensitive group), then he must tell the truth ($Y_i = 0$) & never claims that he belongs to $u_1$.

But if he/she belongs to $u_1$ (ie. the sensitive group), then he/she may tell the truth (i.e $Y_i = 1$) or may tell a lie (ie $Y_i = 0$) or may simply refuse to answer. Each unit in the population has his own way of responding to the question. This model allows different likelihoods for telling the truth, telling a lie or no response for different individuals or units in $u$ ($i = 1, 0$)

Unfortunately, this reasonable assumption on the response complicat the structure of the data collected on the sensitive issue and deprives the statistician from making any meaningful inference on the unknown parameter. To overcome this difficulty the sample units are allowed to keep their privacy by randomizing their responses but the freedom of the individual randomization is replaced by one or two common randomization. Since a randomized response is less revealing then a direct response, it is reasonable to expect greater co-operation.

There are two popular methods, (i) related question method and (ii) unrelated question method.

## Related Question Method

To protect the privacy of the interviwee and to attract his collaboration, we don't measure the sensitive characteristic $Y$, rather we measure a randomized version of it, say $Z$. To do this

each interviwee is presented with two complementary questions:

$Q_0$: Do you belong to $U_0$?

$Q_1$: Do you belong to $U_1$?

Answers to both the questions is either Yes $(Z=1)$ or No $(Z=0)$ and if the reply to $Q_0$ is Yes, then it would be No to $Q_1$, and vice-versa. In deciding which question to answer, each interviwee is provided with an identical randomization device, e.g. a spinner with its face divided into two mutually exclusive and exhaustive parts labelled as $U_0$ & $U_1$. Unseen by the interviewer, each interviwee spins the spinner. He is told to answer $Q_0$ and $Q_1$ depending on whether the spinner stops at $U_0$ and $U_1$. Therefore the outcome of each interviwee is either 1 or 0 and of course without the interviwewer knowing which question has been answered. To be able to utilize the corrected data, it is necessary to assume that the answers are truthful. This means that if the random selection resulted in $Q_0$, then

$$P(Z_i = 1 \mid i \in U_0) = 1 \qquad P(Z_i = 1 \mid i \in U_1) = 0$$

$$P(Z_i = 0 \mid i \in U_0) = 0 \qquad P(Z_i = 0 \mid i \in U_1) = 1$$

On the other hand if the random selection resulted in $Q_1$,

$$P(Z_i = 1 \mid i \in U_1) = 1 \qquad P(Z_i = 1 \mid i \in U_0) = 0$$

$$P(Z_i = 0 \mid i \in U_1) = 0 \qquad P(Z_i = 0 \mid i \in U_0) = 1$$

The above are assumed to be satisfied for every population unit $i$ interviewed with corresponding $Z_i$.

## Estimation of $\phi$

Let, $\pi$ be the probability that the spinner stops at $u_1$, i.e. $Q_1$ is selected for answering. The choice of $\pi$ is under the control of the interviewer and its value is assumed to be known to the data analyst. Under strict randomization $\pi \neq 1$ or $\pi \neq 0$. Further it is chosen to be different from $0.5$. Then the Bernoulli random variable $Z_i$ associated with the population unit $i$ is characterised by the following randomized response distribution denoted by $P_R$.

$$P_R (Z_i = 1) = P(Q_0 \text{ is selected}) \; P(Z_i = 1 \mid Q_0 \text{ is selected}) +$$
$$P(Q_1 \text{ is selected}) \; P(Z_i = 1 \mid Q_1 \text{ is selected})$$

$$= (1-\pi)(1-\phi_i) + \pi \phi_i \;, \text{ where } \phi_i = \begin{cases} 1 & \text{,if } i \in u_1 \\ 0 & \text{,if } i \in u_0 \end{cases}$$

$$= (2\pi - 1) \phi_i + (1-\pi)$$

$$= \theta_i \quad (\text{say})$$

$$\therefore \; P_R (Z_i = 0) = 1 - \theta_i$$

Hence, $\quad \phi_i = \dfrac{\theta_i - (1-\pi)}{(2\pi - 1)}$

Therefore, if our data allow a response based unbiased estimation of $\theta_i$ by $\hat{\theta_i}$, then a responsed based unbiased estimator of $\phi_i$ based on $\theta_i$ will be

$$\hat{\phi_i} = \frac{\hat{\theta_i} - (1-\pi)}{(2\pi - 1)} \quad , i = 1(1)N$$

Note that, $\quad \phi = \dfrac{1}{N} \displaystyle\sum_{i=1}^{N} \phi_i$

$\therefore \; \phi$ is estimated by $\hat{\phi} = \dfrac{1}{n} \displaystyle\sum_{i=1}^{n} \hat{\phi_i}$

$$\therefore \ \hat{\phi} = \frac{\frac{1}{n}\sum_{i=1}^{n}\hat{\theta_i} - (1-\pi)}{(2\pi-1)}, \text{ when, } \frac{1}{n}\sum_{i=1}^{n}\hat{\theta_i} \text{ is the sample}$$

proportion of 'YES' response.

Note that, under SRSWR, the randomized response variables $z_1, z_2, \ldots, z_n$ are all i.i.d Bernoulli random variables with

$$P_R(z_i=1) = (2\pi-1)\phi + (1-\pi) = \theta \ (\text{say}) \qquad , i = 1(1)n$$

Then, $n_1 = \sum_{i=1}^{n} z_i \sim Bin(n, \theta)$

giving $\hat{\theta} = \dfrac{n_1}{n}$

Accordingly, $\hat{\phi} = \dfrac{\hat{\theta} - (1-\pi)}{(2\pi-1)} = \dfrac{\frac{n_1}{n} - (1-\pi)}{(2\pi-1)}$

with $E(\hat{\phi}) = \dfrac{\theta - (1-\pi)}{(2\pi-1)} = \phi$

& $V(\hat{\phi}) = \dfrac{\theta(1-\theta)}{n(2\pi-1)^2}$

The variance is estimated by $\widehat{V(\hat{\phi})} = \dfrac{\hat{\theta}(1-\hat{\theta})}{n(2\pi-1)^2}$

$$= \dfrac{\frac{n_1}{n}\left(1 - \frac{n_1}{n}\right)}{n(2\pi-1)^2}$$

Remark:

Note that, $Var(\hat{\phi}) = \dfrac{\theta(1-\theta)}{n(2\pi-1)^2}$

$$= \dfrac{\{(1-\pi) + (2\pi-1)\phi\}\{\pi - (2\pi-1)\phi\}}{n(2\pi-1)^2}$$

$$= \frac{\pi(1-\pi) + (2\pi-1)\phi\{\pi - 1 + \pi\} - (2\pi-1)^2\phi^2}{n(2\pi-1)^2}$$

$$= \frac{\phi(1-\phi)}{n} + \frac{\pi(1-\pi)}{n(2\pi-1)^2}$$

Therefore, Var $(\hat{\phi})$ is symmetric about $\pi = 0.5$ and this is minimized by taking $\pi$ as far from $0.5$ as possible. To minimize the suspicion of the respondents the value of $\pi$ should not be close to $0$ or $1$, otherwise the respondent might get a strong feeling of being identified as the members belonging to the sensitive group in case the response is 'yes' for large $\pi$ & 'no' for small $\pi$. From this point of view, choice of $\pi$ close to $0.5$ would be most most acceptable but this produces a very large variance of the estimate $\hat{\phi}$. This suggests that the choice of $\pi$ in the interval $[0.6, 0.7]$ or $[0.2, 0.3]$ would be most reasonable.

# Unrelated Question Method

In the related question method, the proportion under study is divided into two mutually exclusive group $U_1$ and $U_0$. One with and the other without the stigmatising version of the ~~sensible~~ sensitive characteristic and based on this division of the population each respondent is asked to respond 'yes' or 'no' depending on whether or not he belongs to the group indicated by one of the two complementary questions $Q_0$ and $Q_1$ which is selected for him to answer. To increase the respondance likelihood of cooperation the related question procedure has been revised. by replacing $Q_0$ ~~and~~ with a harmless dichotomous question as described.

The required form is called the Unrelated Question Procedure in which the population is ~~de~~ divided into two overlapping groups $U_1$ and $U_2$. The group $U_1$ contains all units with the stigmatising version of the sensitive characteristic and the group $U_2$ contains all units having a version of some other dichotomous characteristic. which is completely innocuous and unrelated to the sensitive characteristic under study. Therefore in this procedure $U_1 \cup U_2 \neq U$ and for an obvious reason. $U_2 \subset U_1$ is not allowed.

For example, $Q_1$ and $Q_2$ can be the following

$Q_1$: Do you use illegal drugs?
$Q_2$: Do you prefer football?

It seems reasonable to expect more cooperation from the sample units under this procedure in comparison with the related question procedure.

Under truthful reporting for both procedures the unrelated question procedure can be made to provide a more efficient estimator of the proportion of units $\phi$ associated with the stigmatising version of the sensitive characteristic.

Let, $N_1$ and $N_2$ be the ~~no~~ number of units in $U_1$ and $U_2$ respectively. As before we are interested in estimating $\phi = \frac{N_1}{N}$. We consider the cases separately $\delta$ (known) and $\delta$ (unknown) where $\delta = \frac{N_2}{N}$ (Proportion of individual having some version of the innocuous character).

## Case-1 : $\delta$ known

In this case we select $n$ units by a sampling design of our choice, Each interviewee in provided with the following two questions :

$Q_1$ : Do you belong to $u_1$ ?

$Q_2$ : Do you belong to $u_2$ ?

Unobserved by the interviewer each respondent selects randomly $Q_1$ and/or $Q_2$ with known probabilities $q$ and $(1-q)$ respectively.

The interviewer provides the randomization device. Each respondents answers in either 'YES' or 'NO' depending on which question is randomly selected.

Suppose $n_1$ be the number of 'YES' responses reported, then we have the following results :

Under SRSWR of size $n$ an unbiased estimator of $\phi$ is given by,

$$\hat{\phi} = \frac{n_1}{n} \cdot q + (1-q)\delta \cdot \widehat{P(YES)} = \frac{n_1}{n} = q\hat{\phi} + (1-q)\delta$$

$$\Rightarrow \hat{\phi} = \left\{ \frac{n_1}{n} - (1-q)\delta \right\} / q$$

with $Var(\hat{\phi}) = \dfrac{\alpha(1-\alpha)}{nq^2}$ , where $\alpha = q\phi + (1-q)\delta$ ———(*)

The variance is unbiasedly estimated by :

$$\widehat{Var(\hat{\phi})} = \frac{\hat{\alpha}(1-\hat{\alpha})}{nq^2} \quad ; \quad \hat{\alpha} = \frac{n_1}{n} \quad [\text{proportion of YES response}]$$

Let, $X_i = \begin{cases} 1 & \text{, if ith individual replies 'YES'} \\ 0 & \text{, o.w.} \end{cases}$

So, $X_i \sim Ber\left(p = P(X_i=1) = \alpha\right)$

$\therefore \quad n_1 = \sum\limits_{i=1}^{n} X_i \sim Bin(n, \alpha)$

$E(n_1) = n\alpha \Rightarrow E\left(\frac{n_1}{n}\right) = \alpha$

$Var(n_1) = n\alpha(1-\alpha) \Rightarrow Var\left(\frac{n_1}{n}\right) = \dfrac{\alpha(1-\alpha)}{n}$

$\therefore \quad E(\hat{\phi}) = \dfrac{E\left(\frac{n_1}{n}\right) - (1-q)\delta}{q} = \dfrac{\alpha - (1-q)\delta}{q} = \dfrac{q\phi + (1-q)\delta - (1-q)\delta}{q}$

$E(\hat{\phi}) = \phi$

$$Var(\hat{\phi}) = \frac{Var\left(\frac{m}{n}\right)}{q^2} = \frac{\alpha(1-\alpha)}{nq^2}$$

$$= \frac{\{q\phi + (1-q)\delta\}\{1 - q\phi - (1-q)\delta\}}{nq^2}$$

$$= \frac{1}{nq^2}\left[\{q\phi + (1-q)\delta\}\{1 + q(1-\phi) - (1-q)\delta - q\}\right]$$

$$= \frac{1}{nq^2}\left\{q\phi + (1-q)\delta\right\}\left\{q(1-\phi) + (1-q)(1-\delta)\right\}$$

$$= \frac{q^2\phi(1-\phi)}{nq^2} + \frac{(1-q)^2\delta(1-\delta)}{nq^2} + \frac{q\phi(1-q)(1-\delta) + q(1-\phi)\delta(1-q)}{nq^2}$$

$$= \frac{\phi(1-\phi)}{n} + \frac{(1-q)^2}{nq^2}\delta(1-\delta) + \frac{(1-q)}{nq}\left\{\phi(1-\delta) + \delta(1-\phi)\right\}$$

**Case : 2 :** $\delta$ unknown but $\hat{\delta}$ is available with known variance estimator $Var(\hat{\delta})$

The data collection procedure is same, here we have the following results,

$$\hat{\phi} = \frac{\frac{n_1}{n} - (1-q)\hat{\delta}}{q}$$

$$Var(\hat{\phi}) = \frac{\frac{\alpha(1-\alpha)}{n} + (1-q)^2 Var(\hat{\delta})}{q^2} \quad \left[\begin{array}{l}\text{ignoring the covariance between} \\ \hat{\delta} \text{ and } \frac{n_1}{n}\end{array}\right]$$

$$= \frac{\alpha(1-\alpha)}{nq^2} + \frac{(1-q)^2}{q^2}Var(\hat{\delta})$$

and $\quad \widehat{Var(\hat{\phi})} = \frac{\hat{\alpha}(1-\hat{\alpha})}{nq^2} + \frac{(1-q)^2}{q^2}\widehat{Var(\hat{\delta})}$

**Case 3:** $\delta$ is unknown and no prior estimate of it is available

___

Here if the same randomization device is followed in connecting data from $n$ sampled units, the expression (*) will continue

$$\hat{\phi} = \frac{\frac{n_1}{n} - (1-q)\delta}{q}$$

with $Var(\hat{\phi}) = \frac{\alpha(1-\alpha)}{nq^2}$ where $\alpha = q\phi + (1-q)\delta$ —— (*)

Having an unbiased estimator of $\alpha$ in the usual manner will not permit us to extract exactly one unbiased estimator of $\phi$ from the RHS of (*). Since we have one equations and two unknowns. So the idea is to generate two consistent and independent linear equations of the type. This can be achieved as follows :

We split the total sample size $n$ into two parts $n_1$ and $n_2$. next we perform the randomised response survey on a group $n_1$ sample units selected according to SRSWR using probabilities $q_1$ and $(1-q_1)$ of selecting the questions $B_1$ and $B_2$ respectively. We carry out as similar survey on another group of $n_2$ sample units, selected independently from the entire population again using SRSWR. For this group however the randomization device is made to select the questions $B_1$ and $B_2$ with probabilities $q_2$ and $(1-q_2)$ respectively with $q_1 \neq q_2$.

This randomization procedure results in two equations of the form

$$\alpha_1 = q_1\phi + (1-q_1)\delta$$
$$\alpha_2 = q_2\phi + (1-q_2)\delta$$

let us denote by $n_{11}$ and $n_{21}$, the no. of 'yes' responses for the two groups respectively. Then it is clear that $\frac{n_{11}}{n_1}$ is an u.e. of $\alpha_1$ with variance $\frac{\alpha_1(1-\alpha_1)}{n_1}$

Similarly, $\frac{n_{21}}{n_2}$ is an u.e. of $\alpha_2$, with variance $\frac{\alpha_2(1-\alpha_2)}{n_2}$

∴ To estimate $\phi$. we solve the equations,

$$\frac{n_{11}}{n_1} = q_1\hat{\phi} + (1-q_1)\hat{\delta}$$

$$\frac{n_{21}}{n_2} = q_2 \hat{\phi} + (1-q_2)\hat{\delta}$$

which gives, 

$$q_2 q_1 \hat{\phi} + q_2(1-q_1)\hat{\delta} = q_2 \frac{n_{11}}{n_1} \qquad ——①$$

$$q_1 q_2 \hat{\phi} + q_1 (1-q_2)\hat{\delta} = q_1 \frac{n_{21}}{n_2} \qquad ——②$$

---

① − ②

$$\hat{\delta}\left\{ q_2(1-q_1) - q_1(1-q_2) \right\} = q_2 \frac{n_{11}}{n_1} - q_1 \frac{n_{21}}{n_2}$$

i.e.

$$\hat{\delta} = \frac{q_2 \frac{n_{11}}{n_1} - q_1 \frac{n_{21}}{n_2}}{q_2 - q_1}$$

Again,

$$\hat{\phi} = \frac{\frac{n_{21}}{n_2} - (1-q_2)\hat{\delta}}{q_2} = \frac{\frac{n_{21}}{n_2} - (1-q_2)\left\{ \frac{q_2 \frac{n_{11}}{n_1} - q_1 \frac{n_{21}}{n_2}}{(q_2-q_1)} \right\}}{q_2}$$

$$= \frac{1}{q_2(q_2-q_1)} \left[ \frac{n_{21}}{n_2}q_2 - \frac{n_{11}}{n_1}q_1 - q_2 \frac{n_{11}}{n_1} + q_1 \frac{n_{21}}{n_2} + q_2 \left( q_2 \frac{n_{11}}{n_1} - q_1 \frac{n_{21}}{n_2} \right) \right]$$

$$= \frac{1}{q_2(q_2-q_1)} \left[ \frac{n_{21}}{n_2}q_2 - \frac{n_{21}}{n_2}q_1 - \frac{n_{11}}{n_1}q_2 + \frac{n_{21}}{n_2}q_1 + q_2 \left( q_2 \frac{n_{11}}{n_1} - q_1 \frac{n_{21}}{n_2} \right) \right]$$

$$\hat{\phi} = \frac{(1-q_1)\frac{n_{21}}{n_2} - (1-q_2)\frac{n_{11}}{n_1}}{(q_2-q_1)}$$

$$\hat{\phi} = \frac{(1-q_2)\frac{n_{11}}{n_1} - (1-q_1)\frac{n_{21}}{n_2}}{q_1 - q_2}$$

$$Var(\hat{\phi}) = \frac{1}{(q_1-q_2)^2}\left[ (1-q_2)^2 \frac{\alpha_1(1-\alpha_1)}{n_1} + (1-q_1)^2 \frac{\alpha_2(1-\alpha_2)}{n_2} \right]$$

For a fix total number of observations, $n = n_1 + n_2$ and for fixed values of $q_1, q_2$ and $\delta$ find the optimum $n_1$ and $n_2$ that minimized the variance, we have $n = n_1 + n_2$, so $n_2 = n - n_1$

Thus, $\text{Var}(\hat{\phi}) = \dfrac{(1-q_2)^2 \dfrac{\alpha_1(1-\alpha_1)}{n_1} + (1-q_1)^2 \dfrac{q_2(1-\alpha_2)}{n-n_1}}{(q_1-q_2)^2}$

# Randomized Response Technique for one polychotomous qualitative characteristic

We assume that the sensitive characteristic has $t$ versions (some of them being stigmatizing in nature) and accordingly the population is partitioned into $t$ mutually exclusive and exhaustive parts $U_1, U_2, \ldots, U_t$. We denote by $\phi_i$ the true proportion (unknown) of individuals in the subpopulation $U_i$. Clearly, $\sum_{i=1}^{t} \phi_i = 1$.

## Related Question Method

We prepare $(t-1)$ randomization devices denoted by $D_1, D_2, \ldots, D_{t-1}$. Under any device, the respondent is supposed to respond any of the questions $Q_1, Q_2, \ldots, Q_t$ where,

$Q_j :-$ Do you belong to $U_j$? $(j = 1(1)t)$

However, under the device $i$, the chance of selecting $Q_j$ is $\pi_{ij}$, where $\pi_{ij}$'s are so chosen that

(i) $\quad 0 \leq \pi_{ij} \leq 1$

(ii) $\quad \sum_{j=1}^{t} \pi_{ij} = 1$ for each $i = 1, 2, \ldots, t-1$

(iii) the matrix $A = \left( \left( \pi_{ij} - \pi_{it} \right) \right)_{(t-1) \times (t-1)}$, $\quad i = 1, 2, \ldots, t-1$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad j = 1, 2, \ldots, t-1$

is ~~orthogonal~~. non-singular.

We select $(t-1)$ independent samples of size $n_1, n_2, \ldots, n_{t-1}$ respectively, each drawn according to SRSWR. Each member of the first set of $n_1$ selected respondents is given the randomization device $D_1$ and requested to report the response of "yes" or "no" depending on whether or not he/she belongs to the category yes or "no" for the randomly chosen question. Similarly the other randomization device are used to collect data from the other groups of respondents.

## Estimation of $\phi_1, \phi_2, \ldots, \phi_t$'s

Clearly, it is enough to estimate $\phi_1, \phi_2, \ldots, \phi_{t-1}$, since

$$\phi_t = 1 - \phi_1 - \phi_2 - \cdots - \phi_{t-1}$$

**Notation:** $n_{1i} =$ The no. of 'yes' responses from the individuals using the randomization device $D_i$

$$\alpha_i = \sum_{j=1}^{t} \pi_{ij} \phi_j$$

$$\ell_{9_i} = \alpha_i - \pi_{it} \quad , i = 1(1)\,t-1$$

$$\underset{t-1 \times 1}{\ell_9} = \begin{pmatrix} \ell_{9_1} \\ \ell_{9_2} \\ \vdots \\ \ell_{9_{t-1}} \end{pmatrix} \qquad , \qquad \underset{t-1 \times 1}{\Phi} = \begin{pmatrix} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_{t-1} \end{pmatrix}$$

## Result

An unbiased estimator of $\Phi$ is given by

$$\hat{\Phi} = A^{-1} \hat{\ell_9}$$

where, $\hat{\ell_9} = \begin{pmatrix} \hat{\ell_{9_1}} \\ \hat{\ell_{9_2}} \\ \vdots \\ \hat{\ell_{9_{t-1}}} \end{pmatrix}$ , $\hat{\ell_{9_i}} = \dfrac{n_{1i}}{n_i} - \pi_{it}$ $\quad \left[ \because \hat{\alpha}_i = \dfrac{n_{1i}}{n_i} \ , i = 1(1)\,t-1 \right]$

## Proof:

$$P\left(\text{Yes Response from device } i\right) = \sum_{j=1}^{t} P\left(\begin{array}{l}j^{th} \text{ question selected} \\ \text{from device } i\end{array}\right)$$
$$P\left(\text{Yes} \mid j^{th} \text{ question selected} \atop \text{from device } i\right)$$

$$= \sum_{j=1}^{t} \pi_{ij} \phi_j$$

$$= \alpha_i$$

$\alpha_i$ is unbiasedly estimated by $\alpha_i = \dfrac{n_{i1}}{n_i}$, $i = 1(1)\,t-1$

We have $(t-1)$ equations

$$\frac{n_{i1}}{n_i} = \sum_{j=1}^{t} \pi_{ij}\,\hat{\phi}_j \qquad \left(i = 1(1)\,t-1\right)$$

or, 
$$\frac{n_{i1}}{n_i} = \sum_{j=1}^{t-1} \pi_{ij}\,\hat{\phi}_j + \pi_{it}\,\hat{\phi}_t$$

$$= \sum_{j=1}^{t-1} \pi_{ij}\,\hat{\phi}_j + \pi_{it}\left(1 - \sum_{j=1}^{t-1}\hat{\phi}_j\right)$$

$$= \sum_{j=1}^{t-1}\left(\pi_{ij} - \pi_{it}\right)\hat{\phi}_j + \pi_{it}$$

or, 
$$\sum_{j=1}^{t-1}\underbrace{\left(\pi_{ij} - \pi_{it}\right)}_{a_{ij}}\hat{\phi}_j = \frac{n_{i1}}{n_i} - \pi_{it} = \hat{\ell}_{y_i}$$

or, 
$$\sum_{j=1}^{t-1} a_{ij}\,\hat{\phi}_j = \hat{\ell}_{y_i} \qquad , i = 1(1)\,t-1$$

The system of equations can be written in terms of matrix notation as

$$A\,\hat{\underset{\sim}{\phi}} = \hat{\underset{\sim}{\ell}_y}$$

or, 
$$\hat{\underset{\sim}{\phi}} = A^{-1}\,\hat{\underset{\sim}{\ell}_y} \quad \left[\text{Since, } \pi_{ij} \text{ are so chosen that } A \text{ is non-singular}\right]$$

### Variance of $\hat{\underset{\sim}{\phi}}$

$$\text{Disp}\left(\hat{\underset{\sim}{\phi}}\right) = A^{-1}\,\text{Disp}\left(\hat{\underset{\sim}{\ell}_y}\right)\left(A^{-1}\right)'$$

$$\hat{\underset{\sim}{\ell}_y}_{\,t-1 \times 1} = \begin{pmatrix} \hat{\ell}_{y_1} \\ \vdots \\ \hat{\ell}_{y_{t-1}} \end{pmatrix} , \qquad \hat{\ell}_{y_i} = \frac{n_{i1}}{n_i} - \pi_{it}$$

$$\text{Var}\left(\hat{\ell}_{y_i}\right) = \text{Var}\left(\frac{n_{i1}}{n_i}\right)$$

$$n_{i1} \sim \text{Bin}\left(n_i,\, \alpha_i\right) \text{ independently}$$

So, $Var\left(\hat{\ell}_{gi}\right) = \dfrac{\alpha_i(1-\alpha_i)}{n_i}$, $i = 1(1)\,t-1$

$\&$ $Cov\left(\hat{\ell}_{gi}, \hat{\ell}_{gj}\right) = 0$  $\forall\, i \neq j = 1(1)\,t-1$

$\therefore$ $Disp\left(\hat{\underset{\sim}{\ell}_g}\right) = Diag\left(\dfrac{\alpha_1(1-\alpha_1)}{n_1}, \dfrac{\alpha_2(1-\alpha_2)}{n_2}, \ldots, \dfrac{\alpha_{t-1}(1-\alpha_{t-1})}{n_{t-1}}\right)$

Hence,

$Disp\left(\hat{\underset{\sim}{\Phi}}\right) = A^{-1} Diag\left(\dfrac{\alpha_1(1-\alpha_1)}{n_1}, \dfrac{\alpha_2(1-\alpha_2)}{n_2}, \ldots, \dfrac{\alpha_{t-1}(1-\alpha_{t-1})}{n_{t-1}}\right)(A^{-1})'$

## Remark :

1. Here we explicitly assume that our randomization devices $D_1, \ldots, D_{t-1}$ result in non-singularity of the square matrix $A$. Indeed there are plenty of $\pi_{ij}$'s satisfying the above requirements, for example, we may choose $\pi_{it} = 0$ for $i = 1(1)t$ and choose the other $\pi_{ij}$'s such that $A$ is a circulant matrix of the form $A = \begin{pmatrix} a_1 & a_2 & a_3 & \ldots & a_{t-1} \\ a_{t-1} & a_1 & a_2 & \ldots & a_{t-2} \\ \vdots & \vdots & \vdots & & \vdots \\ a_2 & a_3 & \ldots & a_{t-1} & a_1 \end{pmatrix}$

Here, $0 \leq a_j \leq 1$ and $\displaystyle\sum_{j=1}^{t-1} a_j = 1$ with atleast of the $a_j$'s being strictly positive.

2. To differentiate among various possible choices of the randomization devices $D_1, \ldots, D_{t-1}$, it is necessary to formulate a criterion for comparing the relative performances of the estimators $\hat{\underset{\sim}{\Phi}}$ based on the various choices. Since, the precision of the estimators depends on the choice of $A$, a rational way of choosing $A$ would have been to minimize the determinant of the covariance matrix of the estimators.

ation as

## Unrelated Question Method with 'yes' or 'no' response

Here we assume that there are $t$ stigmatizing versions of the sensitive characteristic thereby partitioning the population $u$ into $t$ subpopulations $U_1, U_2, ...., U_t$. However, there is one unrelated dichotomous characteristic with unknown incidence rate $\delta$ for one of its versions, say $V$.

## Procedure of estimating the proportions by "Greenbarg's procedure"

Using SRSWR scheme of sampling, we draw $t$ independent samples of sizes $n_1, n_2, ....., n_t$ respectively from $t$ randomization devices $D_1, D_2, ..., D_t$.

◻ Each member of the $i^{th}$ group containing $n_i$ individuals asked to select one of the $t$ alternative questions $Q_1, Q_2, ....., Q_t$ with the respective selection probabilities $\pi_{i1}, \pi_{i2}, ..., \pi_{it}$ $(i = 1(1)t)$, where

$Q_j :-$ Do you belong to $U_j$? $(j = 1(1) \overline{t-1})$

and $Q_t :-$ Do you belong to $V$?

<u>Define</u>

$\phi_j$ : population proportion of individuals belonging to $U_j$   $(j=1(1)t)$

$\delta$ : population proportion of individuals belonging to $V$

$$\underset{\sim}{\Phi} = \begin{pmatrix} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_{t-1} \end{pmatrix}$$

$\therefore$ P (Yes response from $i$th randomization device)

$$= \sum_{j=1}^{t} P(\theta_j \text{ is selected from } i\text{th device})$$

$$P(\text{Yes} \mid \theta_j \text{ is selected})$$

$$= \sum_{j=1}^{t-1} a_{ij} \phi_j + q_{it} \cdot \delta = \sum_{j=1}^{t-1} \pi_{ij} \phi_j + \pi_{it} \delta$$

$$= \alpha_i^* \quad (\text{say})$$

$\alpha_i^*$ is unbiasedly estimated by   $\hat{\alpha}_i^* = \dfrac{n_{ii}}{n_i}$

where $n_{ii}$ is the number of individuals giving 'Yes' response from device $D_i$

So we have $t$ equations with $t$ unknowns

$$\frac{n_{ii}}{n_i} = \sum_{j=1}^{t-1} \pi_{ij} \hat{\phi}_i + \pi_{it} \hat{\delta}$$

In matrix notation,   $A^* \begin{pmatrix} \hat{\underset{\sim}{\Phi}} \\ \hat{\delta} \end{pmatrix} = \hat{\underset{\sim}{\alpha}}^*$

where,   $A^* = ((\pi_{ij}))_{t \times t}$   and   $\hat{\underset{\sim}{\alpha}}^* = \begin{pmatrix} \hat{\alpha}_i^* \\ \hat{\alpha}_2^* \\ \vdots \\ \hat{\alpha}_t^* \end{pmatrix}$,

$\pi_{ij}$'s are so choosen that $A^*$ is non-singular

$$\therefore \begin{pmatrix} \hat{\underset{\sim}{\Phi}} \\ \hat{\delta} \end{pmatrix} = A^{*-1} \hat{\underset{\sim}{\alpha}}^*$$

# Unrelated Question Method with multiple response

This is a procedure due to Bourke applicable in situations where the respondents are willing to participate in a multiple choice survey provided an unrelated question with the same number of multiple choices can be blended with the sensitive characteristic.

Suppose that there are $t$ stigmatizing versions of the sensitive characteristic as well as the unrelated characteristic. Let $U_i$ & $U_i^*$ respectively denote the subpopulations corresponding to the $i^{th}$ component of the sensitive characteristic as well as and the innocuous characteristic $(i = 1(1)t)$. It is assumed that the proportions $\delta_1, \delta_2, \dots, \delta_t$ of $U_1^*, U_2^*, \dots, U_t^*$ are known. It is described to estimate the proportions $\phi_1, \phi_2, \dots, \phi_t$ based on the data collected from an SRSWR sample of size $n$.

We prepare a randomization device using, for example, a deck of cards. A known proportion $\pi$ of the cards are of type 1, while the rest are of type 2. Each respondent is supposed to select one card at random. If the selected card is of type 1, then the respondent has to pick up the sensitive issue and is assumed to provide a truthful response to which category he/she belongs. If the selected card is of type-II, then the respondent has to pick up the unrelated question and provide a similar true response. Thus, for example, a respondent belonging to either of the categories $U_i$ and $U_i^*$, would simply respond with the serial number $i$ for data recording

$$\therefore P(\text{Getting Response } i) = P(\text{Sensitive issue selected}) P(\text{Response is } i \mid \text{Sensitive issue selected}) + P(\text{unrelated issue selected}) P(\text{Response is } i \mid \text{Unrelated issue selected})$$

$$= \pi \phi_i + (1-\pi) \delta_i$$

$$= \alpha_i \text{ (say)}$$

$\alpha_i$ is unbiasedly estimated by $\hat{\alpha}_i = \dfrac{n_i}{n}$ , which is the sample proportion of individuals giving response $i$ ( $i = 1(1)t$ )

Thus we obtain $t$ equations in $t$ unknowns,

$$\pi \hat{\phi}_i + (1-\pi) \hat{\delta}_i = \frac{n_i}{n}$$

$$\Rightarrow \hat{\phi}_i = \frac{\frac{n_i}{n} - (1-\pi)\delta_i}{\pi}$$

$$\Rightarrow Var(\hat{\phi}_i) = \frac{1}{\pi^2} \, Var(\hat{\phi}_i) = \frac{\alpha_i(1-\alpha_i)}{n\pi^2}$$

Since, $(n_1, n_2, \ldots, n_t) \sim Multinomial (n, \alpha_1, \alpha_2, \ldots, \alpha_t)$