

Ramakrishna Mission Residential College, Narendrapur

Department of Statistics

Notes on

Double Sampling:

To estimate *population mean or population total*, we often make use of knowledge on an auxiliary variable x at the sampling and/or estimation stage to achieve a gain in precision, e.g., in stratified sampling (stratification on the basis of x), PPS sampling, Ratio and Regression method of estimation etc.

In case of ratio method of estimation, a biased but efficient estimator for \bar{Y}

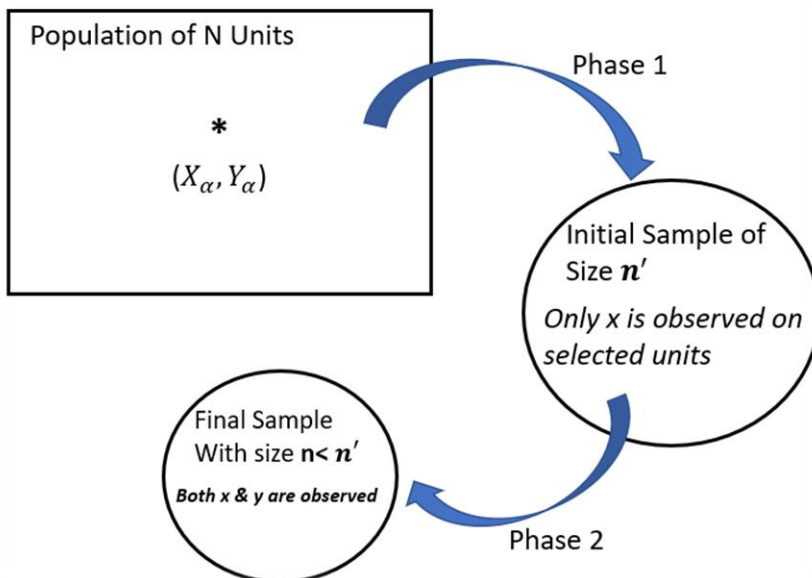
$$\bar{y}_R = \frac{\bar{y}}{\bar{x}} \bar{X} \text{-----(1)}$$

The problem with this estimator is that usually \bar{X} is not known. In that case a double sampling procedure may be used as follows:

- i. From the entire population a sample of size n' is drawn using some suitable sampling scheme and the auxiliary variable x is observed for the selected units. (This information is used to estimate \bar{X} in \bar{y}_R)
- ii. At step two, a sub sample of size n is drawn from the initial sample, or an independent sample of size n is drawn from the entire population using some scheme. Both the study variable and the auxiliary variable are observed for the selected units. Usually, n' is much larger than n .

Note: Since in double sampling, a part of the resource is used to collect information on x from the initial sample, the size of the main sample needs to be reduced as compared to the situation where auxiliary information is not used. However, if the cost of collecting information on x is low, the gain due to incorporating auxiliary information more than offsets the loss due to reduction of the size of the main sample.

The Set up and Notations:



The Population: $U = \{U_1, U_2, \dots, U_\alpha, \dots, U_N\}$

Size of the population: N (Known)

Study variable: y

Auxiliary variable: x (usually cheap to observe and correlated with y)

The Population values: $y(U_\alpha) = Y_\alpha, x(U_\alpha) = X_\alpha, \text{ for } \alpha = 1(1)N$

We can define the parameters on the population values as

Population mean: $\bar{Y} = \sum_1^N Y_\alpha, \bar{X} = \sum_1^N X_\alpha$

Population mean square: $S_y^2 = \frac{1}{N-1} \sum_1^N (Y_\alpha - \bar{Y})^2, S_x^2 = \frac{1}{N-1} \sum_1^N (X_\alpha - \bar{X})^2$

Population covariance: $\sigma_{xy} = \frac{1}{N} \sum_1^N ((Y_\alpha - \bar{Y})(X_\alpha - \bar{X}))$, with $S_{xy} = \frac{N}{N-1} \sigma_{xy}$

1st Phase: From the entire population a sample of size n' (significantly larger than actual sample size n) is drawn using some suitable sampling scheme and the auxiliary variable x is observed for the selected units.

Data: $(x'_i, y'_i), i=1(1)n'$

Only x values are observed.

So, we can calculate $\bar{x}' = \frac{1}{n'} \sum_{i=1}^{n'} x'_i, S_{x'}^2 = \frac{1}{n'-1} \sum_{i=1}^{n'} (x'_i - \bar{x}')^2$

Similarly

we can define (cannot be calculated): $\bar{y}', s_{y'}^2, s'_{x'y'} = \frac{1}{(n'-1)} \sum (x'_i - \bar{x}') (y'_i - \bar{y}')$

2nd Phase: At step two, a sub sample of size n is drawn from the initial sample, or an independent sample of size n is drawn from the entire population using some scheme. Both the study variable and the auxiliary variable are observed for the selected units.

Data: $(x_i, y_i), i=1(1)n$.

In this phase both x and y are observed.

So we observe $\bar{x}, \bar{y}, s_x^2, s_y^2, s_{xy}$ in this phase.

Case 1: A sub sample of size n is drawn from the initial sample and SRSWOR is used at both the phases:

Double sample ratio estimator for population mean is given by

$$\bar{y}_R^d = \frac{\bar{y}}{\bar{x}} \cdot \bar{x}' \text{-----(2)}$$

Let us define $\delta_1 = \frac{\bar{y}-\bar{Y}}{\bar{Y}}, \delta_2 = \frac{\bar{x}-\bar{X}}{\bar{X}}, \delta_3 = \frac{\bar{x}'-\bar{X}}{\bar{X}}$

$$\begin{aligned} \bar{y}_R^d &= (1 + \delta_1)\bar{Y} \left((1 + \delta_2)\bar{X} \right)^{-1} (1 + \delta_3)\bar{X} \\ &= \frac{(1+\delta_1+\delta_3+\delta_1\delta_3)}{(1+\delta_2)} \bar{Y} \end{aligned}$$

$$= (1 + \delta_1 + \delta_3 + \delta_1\delta_3)(1 - \delta_2 + \delta_2^2 - \dots)\bar{Y} \text{ [assuming } |\delta_2| < 1]$$

$$= (1 + \delta_1 + \delta_3 + \delta_1\delta_3 - \delta_2 - \delta_1\delta_2 - \delta_2\delta_3 - \delta_1\delta_2\delta_3 + \delta_2^2 + \dots)\bar{Y}$$

Note that

since at the 2nd phase SRSWOR is used to draw final sample of size n from initial sample of size n', using results of SRSWOR we can have

$$E_2(\bar{y}) = \bar{y}', V_2(\bar{y}) = \left(\frac{1}{n} - \frac{1}{n'}\right) s_{y'}^2, E_2(s_{y'}^2) = s_{y'}^2 \text{-----(3)}$$

$$E_2(\bar{x}) = \bar{x}', V_2(\bar{x}') = \left(\frac{1}{n} - \frac{1}{n'}\right) s_{x'}^2, E_2(s_{x'}^2) = s_{x'}^2 \text{-----(4)}$$

$$\text{also } cov_2(\bar{x}, \bar{y}) = \left(\frac{1}{n} - \frac{1}{n'}\right) s'_{x'y'} \text{ and } E_2(s'_{xy}) = s'_{x'y'} \text{-----(5)}$$

Again

at the 1st phase SRSWOR is used to draw initial sample of size n' from the population of size N. So using results of SRSWOR we can have

$$E_1(\bar{y}') = \bar{Y}, V_1(\bar{y}') = \left(\frac{1}{n'} - \frac{1}{N}\right) S_y^2, E_1(s_{y'}^2) = S_y^2 \text{-----(6)}$$

$$E_1(\bar{x}') = \bar{X}, V_1(\bar{x}') = \left(\frac{1}{n'} - \frac{1}{N}\right) S_x^2, E_1(s_{x'}^2) = S_x^2 \text{-----(7)}$$

$$\text{also } cov_1(\bar{x}', \bar{y}') = \left(\frac{1}{n'} - \frac{1}{N}\right) S_{xy} \text{ and } E_1(s'_{x'y'}) = S_{xy} \text{-----(8)}$$

$$\text{Now, } E(\delta_1) = E_1 E_2 \left(\frac{\bar{y} - \bar{Y}}{\bar{Y}} \right) = \frac{1}{\bar{Y}} \left(E_1(\bar{y}' - \bar{Y}) \right) = 0.$$

Similarly, $E(\delta_2) = 0$.

$$E(\delta_3) = E_1 \left(\frac{\bar{x}' - \bar{X}}{\bar{X}} \right) = 0.$$

Thus, up to first order of approximation, \bar{y}_R^d is unbiased for \bar{Y} .

$$\begin{aligned} \text{Now, } E(\delta_1 \delta_2) &= E \left(\frac{(\bar{y} - \bar{Y})(\bar{x} - \bar{X})}{\bar{X}\bar{Y}} \right) \\ &= \frac{1}{\bar{X}\bar{Y}} \cdot E_1 [E_2 \{ (\bar{x} - \bar{x}' + \bar{x}' - \bar{X})(\bar{y} - \bar{y}' + \bar{y}' - \bar{Y}) \}] \\ &= \frac{1}{\bar{X}\bar{Y}} \cdot E_1 [E_2(\bar{x} - \bar{x}')(\bar{y} - \bar{y}') + (\bar{x}' - \bar{X})E_2(\bar{y} - \bar{y}') + (\bar{y}' - \bar{Y})E_2(\bar{x} - \bar{x}') + (\bar{x}' - \bar{X})(\bar{y}' - \bar{Y})] \\ &= \frac{1}{\bar{X}\bar{Y}} E_1 [cov_2(\bar{x}, \bar{y}) + (\bar{x}' - \bar{X})(\bar{y}' - \bar{Y})] \end{aligned}$$

[since with respect to E_2 (when we are at 2nd phase, i.e., initial sample is treated as population and final sample as sample), $(\bar{x}' - \bar{X})$ and $(\bar{y}' - \bar{Y})$ are constants and $E_2(\bar{y} - \bar{y}') = E_2(\bar{x} - \bar{x}') = 0$]

$$\begin{aligned} &= \frac{1}{\bar{X}\bar{Y}} E_1 \left[\left(\frac{1}{n} - \frac{1}{n'} \right) s'_{x'y'} + (\bar{x}' - \bar{X})(\bar{y}' - \bar{Y}) \right] \quad [\text{by (5)}] \\ &= \frac{1}{\bar{X}\bar{Y}} \left[\left(\frac{1}{n} - \frac{1}{n'} \right) E_1(s'_{x'y'}) + cov_1(\bar{x}', \bar{y}') \right] \\ &= \frac{1}{\bar{X}\bar{Y}} \left[\left(\frac{1}{n} - \frac{1}{n'} \right) S_{xy} + \left(\frac{1}{n'} - \frac{1}{N} \right) S_{xy} \right] \quad [\text{by (8)}] \\ &= \frac{1}{\bar{X}\bar{Y}} \left[\left(\frac{1}{n} - \frac{1}{N} \right) S_{xy} \right] \text{-----(9)} \end{aligned}$$

$$E(\delta_1 \delta_3) = \frac{E_1 E_2 (\bar{y} - \bar{Y})(\bar{x}' - \bar{X})}{\bar{X}\bar{Y}} = \frac{1}{\bar{X}\bar{Y}} [E_1(\bar{x}' - \bar{X})E_2(\bar{y} - \bar{Y})]$$

[with respect to E_2 , $(\bar{x}' - \bar{X})$ is constant]

$$\begin{aligned} &= \frac{1}{\bar{X}\bar{Y}} [E_1 \{ (\bar{x}' - \bar{X})(\bar{y}' - \bar{Y}) \}] \quad [E_2(\bar{y}) = \bar{y}' \text{ by (3)}] \\ &= cov_1(\bar{x}', \bar{y}') = \left(\frac{1}{n'} - \frac{1}{N} \right) \frac{S_{xy}}{\bar{X}\bar{Y}} \quad [\text{by (8)}] \text{-----(10)} \end{aligned}$$

$$E(\delta_2 \delta_3) = \frac{1}{\bar{X}^2} \left[E_1 \left\{ (\bar{x}' - \bar{X}) \left(E_2(\bar{x} - \bar{X}) \right) \right\} \right] = \frac{1}{\bar{X}^2} \left[E_1(\bar{x}' - \bar{X})^2 \right] \quad [E_2(\bar{x}) = \bar{x}' \text{ by (4)}]$$

$$= \frac{V_1(\bar{x}')}{\bar{X}^2} = \left(\frac{1}{n'} - \frac{1}{N}\right) \frac{S_x^2}{\bar{X}^2} \quad [by (7)] \text{-----}(11)$$

$$\begin{aligned} E(\delta_1^2) &= \frac{1}{\bar{Y}^2} E_1 \left(E_2(\bar{y} - \bar{Y})^2 \right) = \frac{1}{\bar{Y}^2} E_1 \left(E_2(\bar{y} - \bar{y}' + \bar{y}' - \bar{Y})^2 \right) \\ &= \frac{1}{\bar{Y}^2} E_1 [V_2(\bar{y}) + (\bar{y}' - \bar{Y})^2 + 2(\bar{y}' - \bar{Y})E_2(\bar{y} - \bar{y}')] \quad [expanding and taking E_2] \end{aligned}$$

[reasons are analogous to those explained earlier]

$$= \frac{1}{\bar{Y}^2} E_1 \left[\left(\frac{1}{n} - \frac{1}{n'}\right) s_{y'}^2 + (\bar{y}' - \bar{Y})^2 \right] \quad [by (3)]$$

$$= \frac{1}{\bar{Y}^2} \left[\left(\frac{1}{n} - \frac{1}{n'}\right) s_y^2 + V_1(\bar{y}') \right] \quad [by (6)]$$

$$= \frac{1}{\bar{Y}^2} \left[\left(\frac{1}{n} - \frac{1}{n'}\right) s_y^2 + \left(\frac{1}{n'} - \frac{1}{N}\right) s_y^2 \right] \quad [by (6)]$$

$$= \frac{1}{\bar{Y}^2} \left(\frac{1}{n} - \frac{1}{N}\right) s_y^2 \text{-----}(12)$$

[Alternatively, $E(\delta_1^2) = \frac{1}{\bar{Y}^2} E(\bar{y} - \bar{Y})^2 = \frac{1}{\bar{Y}^2} [E_1 V_2(\bar{y}) + V_2 E_1(\bar{y})]$]

Similarly, $E(\delta_2^2) = \frac{1}{\bar{X}^2} \left(\frac{1}{n} - \frac{1}{N}\right) s_x^2 \text{-----}(13)$

And $E(\delta_3^2) = \frac{1}{\bar{X}^2} E_1(\bar{x}' - \bar{X})^2 = \frac{1}{\bar{X}^2} V_1(\bar{x}') = \frac{1}{\bar{X}^2} \left(\frac{1}{n'} - \frac{1}{N}\right) s_x^2 \text{-----}(14)$

$$MSE(\bar{y}_R^d) = E(\bar{y}_R^d - \bar{Y})^2 \cong E(\delta_1 + \delta_3 - \delta_2)^2 \bar{Y}^2$$

[Neglecting the terms involving δ_1, δ_2 and δ_3 with power more than 2]

$$\begin{aligned} &= \bar{Y}^2 [E(\delta_1^2 + \delta_2^2 + \delta_3^2 + 2\delta_1\delta_3 - 2\delta_1\delta_2 - 2\delta_2\delta_3)] \\ &= \bar{Y}^2 \left[\left(\frac{1}{n} - \frac{1}{N}\right) \left(\frac{s_y^2}{\bar{Y}^2} + \frac{s_x^2}{\bar{X}^2} - \frac{2S_{xy}}{\bar{X}\bar{Y}}\right) + \left(\frac{1}{n'} - \frac{1}{N}\right) \left(\frac{s_x^2}{\bar{X}^2} + \frac{2S_{xy}}{\bar{X}\bar{Y}} - \frac{2s_x^2}{\bar{X}^2}\right) \right] \\ &= \left(\frac{1}{n} - \frac{1}{N}\right) (S_y^2 + S_x^2 R^2 - 2S_{xy}R) + \left(\frac{1}{n'} - \frac{1}{N}\right) (2RS_{xy} - R^2 S_x^2) \text{-----}(15) \end{aligned}$$

[using (9) to (14)]

A biased but useful estimator of MSE may be obtained by replacing S_y^2, S_x^2, S_{xy}, R by $s_y'^2, s_x'^2, s_{xy}'$ and $\frac{\bar{y}}{\bar{x}}$.

Case 2: The main sample is drawn from the entire population, from independent of the initial sample.

In that case, $E(\delta_1\delta_3)=0$

$$E(\delta_2\delta_3)=0$$

$$E(\delta_1\delta_2) = = \frac{cov(\bar{x}, \bar{y})}{\bar{X}\bar{Y}} = \left(\frac{1}{n} - \frac{1}{N}\right) \frac{S_{xy}}{\bar{X}\bar{Y}}$$

$$E(\delta_1^2) = \frac{s_y^2}{Y^2} \left(\frac{1}{n} - \frac{1}{N} \right)$$

$$E(\delta_2^2) = \frac{s_x^2}{X^2} \left(\frac{1}{n} - \frac{1}{N} \right)$$

$$E(\delta_3^2) = \frac{1}{X^2} \left(\frac{1}{n'} - \frac{1}{N} \right) S_x^2$$

[using the results of SRSWOR]

$$\text{MSE}(\bar{y}_R^d) = \left(\frac{1}{n} - \frac{1}{N} \right) [S_y^2 + R^2 S_x^2 - 2R S_{xy}] + \left(\frac{1}{n'} - \frac{1}{N} \right) R^2 S_x^2 \text{-----(16)}$$

Double Sampling Regression Estimator:

Usual regression estimator is given by $\bar{y}_{Reg}' = \bar{y} - b(\bar{x} - \bar{X})$.

When \bar{X} is not known, it can be replaced by \bar{x}' , the sample mean obtained from the initial sample of size n' , say.

Thus, we have the double sampling regression estimator as $\bar{y}_{Reg}^d = \bar{y} - b(\bar{x} - \bar{x}')$

$$\text{Let } \delta_1 = \frac{\bar{y} - \bar{Y}}{\bar{Y}}, \delta_2 = \frac{\bar{x} - \bar{X}}{\bar{X}}, \delta_3 = \frac{\bar{x}' - \bar{X}}{\bar{X}}, \delta_4 = \frac{b - \beta}{\beta}$$

$$\text{Then, } \bar{y}_{Reg}^d = \bar{Y}(1 + \delta_1) + (1 + \delta_4)\beta\bar{X}(\delta_3 - \delta_2)$$

$$\begin{aligned} \text{MSE}(\bar{y}_{Reg}^d) &= E[(\bar{y}_{Reg}^d - \bar{Y})^2] = E[\{\delta_1\bar{Y} + \beta\bar{X}(\delta_3 - \delta_2 + \delta_3\delta_4 - \delta_2\delta_4)\}^2] \\ &= \bar{Y}^2 E(\delta_1^2) + \beta^2 \bar{X}^2 E(\delta_3 - \delta_2 + \delta_3\delta_4 - \delta_2\delta_4)^2 + 2\beta\bar{X}\bar{Y} E(\delta_1\delta_3 - \delta_2\delta_1 + \dots) \\ &\cong \bar{Y}^2 E(\delta_1^2) + \beta^2 \bar{X}^2 E(\delta_3^2 + \delta_2^2 - 2\delta_2\delta_3) + 2\beta\bar{X}\bar{Y} E(\delta_1\delta_3 - \delta_2\delta_1) \end{aligned}$$

[Neglecting the terms in δ_1, δ_2 and δ_3 in power greater than 2]

$$\begin{aligned} &= \left(\frac{1}{n} - \frac{1}{N} \right) S_y^2 + \beta^2 \left(\frac{1}{n'} - \frac{1}{N} \right) S_x^2 + \beta^2 \left(\frac{1}{n} - \frac{1}{N} \right) S_x^2 - 2\beta^2 \left(\frac{1}{n'} - \frac{1}{N} \right) S_x^2 + 2\beta \left(\frac{1}{n'} - \frac{1}{N} \right) S_{xy} - \\ &2\beta \left(\frac{1}{n} - \frac{1}{N} \right) S_{xy} \\ &= \left(\frac{1}{n} - \frac{1}{N} \right) [S_y^2 + \beta^2 S_x^2 - 2\beta S_{xy}] + \left(\frac{1}{n'} - \frac{1}{N} \right) (2\beta S_{xy} - \beta^2 S_x^2) \\ &= \left(\frac{1}{n} - \frac{1}{N} \right) (1 - \rho^2) S_y^2 + \left(\frac{1}{n'} - \frac{1}{N} \right) \rho^2 S_y^2 \end{aligned}$$

Note that if the main sample is drawn independently of the initial sample, we get the same result.

Optimum choice of n and n'

For MSE of both ratio and regression estimator in case of double sampling, we may write

$$\text{MSE} = \frac{V_1}{n} + \frac{V_2}{n'}$$
 with proper choice of V_1 & V_2 .

A simple cost function in case of double sampling may be given as,

$C_0 = c_1 n' + c_2 n$ where c_1 is average cost of surveying one unit in the initial sample and c_2 is that for the final sample.

If we want optimum choice of n and n' by minimizing MSE when cost is fixed at c_0 , say.

Then the objective function can be taken as follows

$$\varphi = \frac{V_1}{n} + \frac{V_2}{n'} + \lambda(c_1 n' + c_2 n - c_0)$$

$$\frac{\partial}{\partial n} \varphi = 0 \Rightarrow \frac{V_1}{n^2} = \lambda c_2 \Rightarrow n = \sqrt{\frac{V_1}{\lambda c_2}}$$

Similarly, $n' = \sqrt{\frac{V_2}{\lambda c_1}}$

$$\text{Also, } c_0 = c_1 n' + c_2 n = \frac{\sqrt{c_1 V_2}}{\sqrt{\lambda}} + \frac{\sqrt{c_2 V_1}}{\sqrt{\lambda}}$$

$$\Rightarrow \frac{1}{\sqrt{\lambda}} = \frac{c_0}{\sqrt{c_1 V_2} + \sqrt{c_2 V_1}}$$

$$\therefore n_{opt} = \frac{C_0}{\sqrt{C_1 V_2} + \sqrt{C_2 V_1}} \sqrt{\frac{V_1}{C_2}}$$

$$\therefore n'_{opt} = \frac{C_0}{\sqrt{C_1 V_2} + \sqrt{C_2 V_1}} \sqrt{\frac{V_2}{C_1}}$$

$$= n_{opt} \sqrt{\frac{V_2}{V_1} \times \frac{C_2}{C_1}}$$

Comparison of Performance of Mean per unit Estimator in Direct Sampling with that of Regression Estimator in case of Double Sampling

In case of double sampling, a sample of size n' , is chosen to observe x , and then a sub sample of size n is drawn from initial sample, and y is observed. If c' and c respectively be the average cost of surveying x and y for one unit. Then for the cost of survey in case of double sampling is $c'n' + cn$. Usually n' is much larger than n . Suppose we draw a sample of size n_0 . Then we must have $n_0 = \frac{c'n'}{c} + n$ -----(17)

so that the cost of survey remains same for both the schemes.

Now MSE for double sampling regression estimator is

$$\left(\frac{1}{n} - \frac{1}{N}\right) (1 - \rho^2) S_y^2 + \left(\frac{1}{n'} - \frac{1}{N}\right) \rho^2 S_y^2 = \left(\frac{1}{n} - \frac{1}{N}\right) S_y^2 - \left(\frac{1}{n} - \frac{1}{n'}\right) \rho^2 S_y^2$$

And variance of mean per unit estimator in case of direct sampling is

$$\left(\frac{1}{n_0} - \frac{1}{N}\right) S_y^2$$

Hence, double sampling regression estimator will be beneficial if,

$$\begin{aligned} &\left(\frac{1}{n} - \frac{1}{N}\right) S_y^2 - \left(\frac{1}{n} - \frac{1}{n'}\right) \rho^2 S_y^2 < \left(\frac{1}{n_0} - \frac{1}{N}\right) S_y^2 \\ \Rightarrow \rho^2 \left(\frac{1}{n'} - \frac{1}{n}\right) &< \left(\frac{1}{n_0} - \frac{1}{n}\right) \\ \Rightarrow \left(\frac{1}{n} - \frac{1}{n'}\right) \rho^2 &> \left(\frac{1}{n} - \frac{1}{n_0}\right) \\ \Rightarrow \rho^2 > \frac{\frac{1}{n} - \frac{1}{n_0}}{\frac{1}{n} - \frac{1}{n'}} &= \left(\frac{n_0 - n}{n_0}\right) \frac{n'}{n' - n} = \frac{\frac{c'n'}{c}}{\frac{c'n'}{c} + n} \cdot \frac{1}{1 - \frac{n}{n'}} = \frac{1}{\left(1 + \frac{cn}{c'n'}\right) \left(1 - \frac{n}{n'}\right)} \end{aligned}$$