



Proportional Stratification Component of Design Effect

Nidal Mohamed Mustafa Abd Elsalam

Department of Statistics & Computation Faculty of Technology of Mathematical Sciences & Statistics

Al Neelain University

*Corresponding Author: Nidal Mohamed Mustafa, [Email: nidalm2 @ gmail.com](mailto:Email:nidalm2@gmail.com)

ABSTRACT

This paper is concerned with Stratification as a component of design effect (Deff). There is a need to determine the pattern and magnitude of this component when the sample allocation is proportional, and to see how such component behaves when parameters like strata means, strata variances, strata proportions and sample size of each stratum vary. Factorial analysis is carried out, to determine the effects of those various factors on the design effect and to form some idea about its expected magnitude when faced with a design of a specific given structure. Also, an attempt will be made to arrive at the expressions, which are believed original, for the design effect in the case of a population consisting of two strata. It is shown that the different combinations of strata means and strata variances affect significantly the design effect. The Deff values increase with the increase in the levels of variances of strata, and as differences in strata means become less. But decrease as strata sizes tend to equality. The precision of proportionate stratified sampling is investigated and examined, and consequently, it is confirmed that it is the same as that of simple random sampling (SRS) when strata means are equal.

المستخلص

تختص هذه الورقة بدراسة المعاينة الطبقيّة كمكون أساسي لأثر التصميم (Deff). وقد استدعي ذلك تحديد نمط و سلوك هذا المكون في حالة التقسيم المتناسب، تحديداً عندما تتغير معالم المجتمع التي تتمثل في الوسط الحسابي والتباين للطبقات، النسبة الطبقيّة و حجم العينة في كل طبقة. و لتحديد اثار تلك العوامل علي أثر التصميم، تم استخدام التحليل العاملي حتي يتم تكوين فكرة عن النمط المتوقع عندما يواجه الباحث بتصميم مماثل. و هناك محاولة ايضا لاشتقاق صيغ لأثر التصميم بدلالة المعالم المؤثرة. و يعتقد أن هذه الصيغ تمثل إضافة أصيلة. و قد تم التوصل إلي أن الأوساط الحسابية و التباينات للطبقات تؤثر بصورة فعالة علي أثر التصميم، مؤدية الي ارتفاع و إنخفاض قيمته مع ارتفاع مستويات تباينات الطبقات و الفروقات في الأوساط الحسابية و إنخفاضه مع تقارب حجم الطبقات، علي التوالي. بالإضافة الي ذلك، تم تحري و فحص دقة المعاينة الطبقيّة في حالة التقسيم

المتناسب و عليه, تم التأكيد علي أنها تساوي دقة المعاينة العشوائية البسيطة عندما يتساوي الوسط الحسابي للطبقات.

KEYWORDS: *Population Design Effect (Deff), proportional stratification, factorial analysis.*

INTRODUCTION

The design effect (Deff) continues to be a valuable tool in sample surveys. It is widely known that it can be defined in the following form:

$$\text{Design Effect} = \frac{\text{Variance of an estimate for a given design}}{\text{Variance of the estimate for a simple random sample of the same size}}$$

Or in another form (as in case of stratified random sampling):

$$Deff = \frac{V(\bar{y}_{st})}{V(\bar{y})}$$

Where; \bar{y}_{st} refers to the stratified mean. And in fact, Deff, plays a very important role to guide every sampler in the design and analysis stage, Kish⁽¹⁾. And it is usually affected by stratification as a complex sample design, kish^(2,3), thus it becomes inevitably essential to explore this effect. And in this context, the paper is confined only to the population mean as the main estimate. This paper extends the previous work of Kalton et al^(4&5) on design effect in proportional stratification. The design effect was expressed only in its general form without referring to the factors affecting its behavior and pattern. Its' value is

believed to range from value 1 to 7, 8 or even up to 30, Curtin L, et al⁽⁶⁾ and David A Lacher⁽⁷⁾. The main objectives of this paper are: to derive appropriate expressions for the design effect in terms of factors expected to affect its pattern, and to combine those factors in a factorial way so as to reach to clear results about the behavior and magnitude of design effect in proportional stratification. The study shows that the different combinations of strata means and strata variances for two strata case only, affect significantly the design effect. The Deff values increase with the increase in the levels of variances of strata, and as differences in means

become less. But decrease as strata sizes tend to equality. While, the precision of proportionate stratified sampling is the same as that of simple random sampling (SRS) when strata means are equal.

MATERIALS AND METHODS

Consider a population of size N that is divided into L non-overlapping strata. Let N_h be the size of stratum h , and suppose a sample of size n_h is drawn by simple random sampling (SRS). It is well known that the variance of a stratified mean, (ignoring the finite population correction factor) Cochran⁽⁸⁾,

is expressed as: $V(\bar{y}_{st}) = \sum_h W_h^2 \frac{S_h^2}{n_h}$

Where; W_h : Proportion of stratum h which is expressed as N_h/N S_h^2 : Variance of stratum h or

$$S_h^2 = \frac{\sum_i (Y_{hi} - \bar{Y}_h)^2}{(N_h - 1)}$$

Assuming that we have only two strata, we consider the behavior of the design effect of the mean at different levels of the factors:

- proportion of stratum 1 (W_1)
- sampling weight (w) which includes sampling weight for each stratum ($w_h = N_h/n_h$)
- strata variances
- strata means

And this will be done under the following conditions:

- Equal strata variances but unequal strata means.
- Equal strata means but unequal strata variances.
- Unequal strata means & unequal strata variances.

Recalling proportional allocation, we substitute:

$$N_h = \frac{nN_h}{N}$$

This yields the following form which represents the variance of stratified mean in proportional sampling as:

$$V_{prop}(\bar{y}_{st}) = \sum_h \frac{N_h^2}{N^2} \frac{S_h^2}{nN_h} N = \frac{\sum_h W_h S_h^2}{n}$$

→(1)

While the variance of the mean for a simple random sample is derived in the following way: Recalling that

$$(N-1)S^2 = \sum_h \sum_i^{N_h} (y_{hi} - \bar{Y})^2$$

$$= \sum_h \sum_i^{N_h} (y_{hi} - \bar{Y} - \bar{Y}_h + \bar{Y}_h)^2$$

$$= \sum_h \sum_i^{N_h} (y_{hi} - \bar{Y}_h)^2 + \sum_h \sum_i^{N_h} (\bar{Y}_h - \bar{Y})^2$$

(Since the sum of the deviations equals zero)

$$= \sum_h^l (N_h - 1)S_h^2 + \sum_h^l N_h (\bar{Y}_h - \bar{Y})^2$$

∴

$$S^2 = \frac{1}{(N-1)} \left[\sum_h^l (N_h - 1)S_h^2 + \sum_h^l N_h (\bar{Y}_h - \bar{Y})^2 \right]$$

Thus, the variance of simple random sampling in case of two strata and ignoring the correction factor can be stated as:-

$$\begin{aligned} V(\bar{y}) &= \frac{S^2}{n} = \frac{1}{(N-1)n} \left[\sum_h^l (N_h - 1)S_h^2 + \sum_h^l N_h (\bar{Y}_h - \bar{Y})^2 \right] \\ &= \\ &= \frac{1}{(N-1)n} \{ [(N_1 - 1)S_1^2 + (N_2 - 1)S_2^2] + [N_1(\bar{Y}_1 - \bar{Y})^2 + N_2(\bar{Y}_2 - \bar{Y})^2] \} \end{aligned} \rightarrow (2)$$

Such that;

$$\bar{Y} = \frac{(N_1 \bar{Y}_1 + N_2 \bar{Y}_2)}{N}$$

And the design effect in this case will be in the following form:

$$\begin{aligned} Deff &= \frac{V_{prop}(\bar{y}_{st})}{V(\bar{y})} \\ &= \\ &= \frac{\sum_h W_h S_h^2}{\frac{1}{(N-1)} [(N_1 - 1)S_1^2 + (N_2 - 1)S_2^2 + N_1(\bar{Y}_1 - \bar{Y})^2 + N_2(\bar{Y}_2 - \bar{Y})^2]} \end{aligned}$$

→(3)

Where;

\bar{Y}_1
 \bar{Y}_2 : are the stratum 1 and stratum 2 mean respectively

\bar{Y} : Population mean

Now moving to the conditions mentioned above we get:

Equal Strata Variances & Unequal Strata Means

Considering the case of equal strata variances and unequal strata means, equation (3) becomes in the following form: $Deff =$

$$\frac{S_c^2(N-1)}{S_c^2(N-2) + [N_1(\bar{Y}_1 - \bar{Y})^2 + N_2(\bar{Y}_2 - \bar{Y})^2]} \rightarrow (4)$$

Where S_c^2 denotes the common variance between the two strata. Using the following combinations (hypothetical values) of means and variances of table (1), together with 10 levels of W_1 (0.05 to 0.50, in steps of 0.05), we get 60 values of Deff represented in table (2). Table (2) shows that there are six design effects combinations due to proportionate sampling in case of equal strata variances, denoted by S_c^2 and different strata means (\bar{Y}_1 & \bar{Y}_2) which were mentioned in table (1). Note that the factor w is not involved, because it does not appear in formula number (4). The values of the design effect, as seen

from the table, are the relative stratum 1 size, the common strata variance and the strata means change in proportionate allocation. Thus, the Deff in this case is always less than 1, although as expected, increases as differences in means become less and decreases as strata sizes tend to equality, Salganik ⁽⁹⁾. The effect

of increase in variance from S_c^2 (100) to S_c^2 (150) increases the design effect for combinations (1.4) to (1.6) of table (2). Table (2) Design effects due to proportionate sampling in case of equal strata variances & different strata means for combinations (1.1) to (1.6) of table 2.

Table (1) Combinations of means & variances for case of Equal Strata Variances & Unequal Strata Means

Combination	S_c^2	\bar{Y}_1	\bar{Y}_2	\bar{Y}_1/\bar{Y}_2
1.1	100	50	12.50	0.25
1.2	100	50	25.0	0.50
1.3	100	50	37.50	0.75
1.4	150	50	12.50	0.25
1.5	150	50	25.0	0.50
1.6	150	50	37.50	0.75

Table (2) Design effects due to proportionate sampling in case of equal strata variances & different strata means for combinations (1.1) to (1.6)

W1	W2	N1	N2	Deff (1.1)	Deff (1.2)	Deff (1.3)	Deff (1.4)	Deff (1.5)	Deff(1.6)
0.05	0.95	250	4750	0.599555	0.771168	0.931070	0.691945	0.834894	0.955160
0.10	0.90	500	4500	0.441369	0.640036	0.876844	0.542382	0.727339	0.921460
0.15	0.85	750	4250	0.358022	0.556534	0.833988	0.455508	0.653101	0.895873
0.20	0.80	1000	4000	0.307669	0.500000	0.800096	0.399984	0.600024	0.876044
0.25	0.75	1250	3750	0.274948	0.460424	0.773499	0.362586	0.561417	0.860146
0.30	0.70	1500	3500	0.252939	0.432421	0.753017	0.336820	0.533340	0.846678
0.35	0.65	1750	3250	0.238115	0.412889	0.737822	0.319179	0.513372	0.834350
0.40	0.60	2000	3000	0.228547	0.399984	0.727339	0.307669	0.500000	0.822024
0.45	0.55	2250	2750	0.223166	0.392621	0.721191	0.301153	0.492306	0.808691
0.50	0.50	2500	2500	0.221429	0.390227	0.719164	0.299041	0.489794	0.793482

Equal Strata Means & Unequal Strata Variances

We now examine the case where the strata variances are taken different, while the two strata means are equal. And from Equation (1):

$$V_{prop}(\bar{y}_{st}) = \frac{\sum_h W_h S_h^2}{n_h}$$

$$= \frac{1}{n} [W_1 S_1^2 + W_2 S_2^2] \text{ With}$$

$$S_h^2 = \frac{1}{(N_h - 1)} \left[\sum_i^{N_{hi}} y_{hi}^2 - N_h (\bar{Y}_h)^2 \right]$$

Implying that

$$S_1^2 = \frac{1}{(N_1 - 1)} \left[\sum_i^{N_{1i}} y_{1i}^2 - N_1 (\bar{Y}_1)^2 \right]$$

$$\text{And } S_2^2 = \frac{1}{(N_2 - 1)} \left[\sum_i^{N_{2i}} y_{2i}^2 - N_2 (\bar{Y}_2)^2 \right] \quad 47$$

And since the two strata means are equal, $\bar{Y}_1 = \bar{Y}_2 = \bar{Y}$;

$$\sum_i^{N_1} y_{1i}^2 = (N_1 - 1) S_1^2 + N_1 (\bar{Y})^2$$

And ;

$$\sum_i^{N_2} y_{2i}^2 = (N_2 - 1) S_2^2 + N_2 (\bar{Y})^2 \text{ And}$$

thus the design effect is expressed in the form:

$$Deff = \frac{\{W_1 \sum_i^{N_1} (y_{1i}^2 - N_1 \bar{Y}^2) / (N_1 - 1)\} + \{W_2 \sum_i^{N_2} (y_{2i}^2 - N_2 \bar{Y}^2) / (N_2 - 1)\}}{\frac{1}{(N-1)} [(N_1 - 1) S_1^2 + (N_2 - 1) S_2^2]}$$

→ (5)

Equation (5) is applied using the different combinations of table (3)

Table (3) :Mean-variance combinations for case Equal Strata Means & Unequal Strata Variances

Combination	S ₂ 1	S ₂ 2	\bar{Y}
2.1	15	20	20
2.2	20	100	20
2.3	600	700	20
2.4	15	20	150
2.5	20	100	150
2.6	600	700	150

resulting in table (4) which contains the design effects 1 to 6, of equal strata means & different strata variances From table (4), the Deff values are similar in the lower level of \bar{Y} to those of the high level of it indicating that using the same level of S_h^2 , regardless of any level of \bar{Y} , while the two strata means are equal, does not affect the values of the design effect. But all the values of table (4) are, up to three decimal places, equal and are equal to 1. This means that precision of proportionate stratified sampling is the same as that of SRS when strata means are equal. And this conforms to the result in sampling theory, Cochran⁽⁸⁾.

Table (4) Design effects due to proportionate sampling in case of equal strata means & different strata variances for combinations (2.1) to (2.6)

W1	W2	N1	N2	Deff (2.1)	Deff (2.2)	Deff (2.3)	Deff (2.4)	Deff (2.5)	Deff(2.6)
0.05	0.95	250	4750	1.00015	1.00005	1.00017	1.00015	1.00005	1.00017
0.10	0.90	500	4500	1.00016	1.00006	1.00018	1.00016	1.00006	1.00018
0.15	0.85	750	4250	1.00016	1.00007	1.00018	1.00016	1.00007	1.00018
0.20	0.80	1000	4000	1.00017	1.00009	1.00018	1.00017	1.00009	1.00018
0.25	0.75	1250	3750	1.00017	1.00010	1.00019	1.00017	1.00010	1.00019
0.30	0.70	1500	3500	1.00018	1.00012	1.00019	1.00018	1.00012	1.00019
0.35	0.65	1750	3250	1.00018	1.00013	1.00019	1.00018	1.00013	1.00019
0.40	0.60	2000	3000	1.00019	1.00015	1.00019	1.00019	1.00015	1.00019
0.45	0.55	2250	2750	1.00019	1.00018	1.00020	1.00019	1.00018	1.00020
0.50	0.50	2500	2500	1.00020	1.00020	1.00020	1.00020	1.00020	1.00020

48

Unequal Strata Means & Unequal Strata Variances In this section, different strata means, as well as

different strata variances are used. The design effect form, accompanied with these combinations is:

$$Deff = \frac{\{[W_1 \sum_i^{N_1} (y_{1i}^2 - N_1 \bar{Y}_1^2)]/(N_1 - 1)\} + \{[W_2 \sum_i^{N_2} (y_{2i}^2 - N_2 \bar{Y}_2^2)]/(N_2 - 1)\}}{\frac{1}{(N - 1)} \{[(N_1 - 1)S_1^2 + (N_2 - 1)S_2^2] + [N_1(\bar{Y}_1 - \bar{Y})^2 + N_2(\bar{Y}_2 - \bar{Y})^2]\}} \rightarrow (6)$$

The mean-variance combinations here are taken from table (5), and according

Table (5): Combinations of means & variances for case Unequal Strata Means & Unequal Strata Variances

Combinations	s^2	s^2	\bar{Y}_1	\bar{Y}_2
3.1	15	20	20	9
3.2	20	100	20	9
3.3	600	700	20	9
3.4	15	20	150	100
3.5	20	100	150	100
3.6	600	700	150	100

Table (6) is formed, which is about the design effect values resulting from using combinations from 1 to 6. It can be realized that in columns (3.1) to (3.3), the Deff values are high in this part than that of the columns (3.4) to (3.6). This is due to the low levels of $\bar{Y}_1(20)$ & $\bar{Y}_2(9)$, and the values increase with the increase in the levels of variances of strata, but decrease as strata sizes tend to equality. And also, the general pattern of its relation to differences in means is retained

Factorial Model

The factorial model, Montgomery⁽¹⁰⁾ and Hinkelmann⁽¹¹⁾ et al, which represents the combination of the proportion of stratum h (W_i), the three cases (CASE) of the combinations of proportional allocation (equal strata variances but unequal strata means, equal strata means but unequal strata variances, unequal strata means & unequal strata variances), and MVC which denotes the mean variance combination, is shown as follows:

$$Y_{ijk} = \mu + W_i + \text{CASE}_j + \text{MVC}_k + (W_i \text{CASE})_{ij} + (W_i \text{MVC})_{ik} + (\text{CASEMVC})_{jk} + (W_i \text{CASEMVC})_{ijk} + e_{ijk}$$

Where; $i=1, 2, 3, \dots, 10$

$j=1, 2, 3.$

$k=1, 2, \dots, 6$

Table (7) shows the result of the factorial model considered here. It is clearly seen that all the main effects (W_i , CASE, and MVC) and two-factor interactions (except for $W_i \times MVC$) are significantly affecting the design effect.

Table (6) Design effects due to proportionate sampling in case of different strata means & different strata variances for combinations (3.1) to (3.6)

W1	W2	N1	N2	Deff (3.1)	Deff (3.2)	Deff (3.3)	Deff (3.4)	Deff (3.5)	Deff(3.6)
0.05	0.95	250	4750	0.774643	0.943546	0.991968	0.142578	0.446992	0.854173
0.10	0.90	500	4500	0.641678	0.894189	0.984631	0.079741	0.290185	0.754162
0.15	0.85	750	4250	0.555116	0.850865	0.978142	0.056942	0.216319	0.682481
0.20	0.80	1000	4000	0.495299	0.812720	0.972485	0.045338	0.173528	0.629655
0.25	0.75	1250	3750	0.452475	0.779089	0.967649	0.038454	0.145763	0.590180
0.30	0.70	1500	3500	0.421299	0.749460	0.963628	0.034032	0.126436	0.560679
0.35	0.65	1750	3250	0.398649	0.723448	0.960419	0.031084	0.112350	0.539013
0.40	0.60	2000	3000	0.382633	0.700775	0.958024	0.029121	0.101780	0.523813
0.45	0.55	2250	2750	0.372117	0.681269	0.956449	0.027882	0.093723	0.514232
0.50	0.50	2500	2500	0.366473	0.664864	0.955705	0.027232	0.087577	0.509806

Table (7): General Linear Model, Case iv: Deff versus W1, CASE, and MVC in proportionate sampling for combinations of Means & Variances Analysis of Variance for Deff, using Adjusted SS for Tests Factor Type Levels Values

W1 fixed 10 1 2 3 4 5 6 7 8 9 10
 CASE fixed 3 1 2 3
 MVC fixed 6 1 2 3 4 5 6

Source	DF	Seq SS	Adj SS	Adj MS	F	P
W1	9	0.39863	0.39863	0.04429	22.77	0.000
CASE	2	11.33699	11.33699	5.66850	2914.14	0.000
MVC	5	4.35495	4.35495	0.87099	447.77	0.000
W1*CASE	18	0.20858	0.20858	0.01159	5.96	0.000
W1*MVC	45	0.07646	0.07646	0.00170	0.87	0.688
CASE*MVC	10	6.90918	6.90918	0.69092	355.20	0.000
Error	90	0.17507	0.17507	0.00195		
Total	179	23.45986				

Table (8) shows some of the details of coefficients, standard errors (S.E), T-value & P- values for the various effects and interactions, respectively. Levels 1, 2, 3 and 4 of W_1 , for example, affect the

Deff, averaged over all levels of the other two factors, by increasing it by: 0.090951, 0.042245, 0.016538 and 0.000045 respectively. And the interpretation will move on like the same way.

Table (8): Coefficients & P-value for Main Effects

Term	Coef	SE Coef	T	P
Constant	0.648419	0.003287	197.25	0.000
1	0.106605	0.009862	10.81	0.000
2	0.056289	0.009862	5.71	0.000
3	0.025360	0.009862	2.57	0.012
4	0.003948	0.009862	0.40	0.690
5	0.011667	0.009862	1.18	0.240
6	-0.023295	0.009862	-2.36	0.020
7	-0.031938	0.009862	-3.24	0.002
8	-0.038185	0.009862	-3.87	0.000
9	-0.042386	0.009862	-4.30	0.000
CASE				
1	-0.216905	0.004649	-46.66	0.000
2	0.351741	0.004649	75.66	0.000
MVC				
1	-0.048155	0.007351	6.55	0.000
2	0.110171	0.007351	14.99	0.000
3	0.270415	0.007351	36.79	0.000
4	-0.029924	0.007351	-4.07	0.000
5	-0.225193	0.007351	-30.64	0.000

RESULTS:

The study showed that the derived expressions for the design effect in case of stratification with proportional allocation helped very much to reveal the impact of some factors on the pattern of the Deff. And the results which were reached through the factorial analysis are of great value to any sampler. For more elaboration, the basic findings of the paper are as follows: *The values of the design effect in proportionate sampling in case of equal strata variances and different strata means are always less than "1". And as expected, increase as differences in means become less and decrease as strata sizes tend to equality. Also, the effect of increase in variance from a certain level to a higher one increases the level of design effect.*Regarding the case of equal strata means and different strata variances, it is clearly seen that using the same level of S_h^2 , regardless of any level of \bar{Y} , while the two strata means are equal, does not affect the values of the design effect. But all the values are equal and are equal to "1", up to three decimal places,. This means that precision of proportionate stratified sampling is the same as that of SRS when strata means are equal. *For the case of different strata means and different strata variances, the Deff values increase with the increase in the levels of variances of strata, but decrease as strata sizes tend to equality. And also, increase as differences in means become less.

CONCLUSION:

An attempt has been made to find appropriate and relevant expressions for the design effect in proportional stratification sampling. This is through taking into consideration the combinations of the various circumstances affecting its behavior. It is hoped that the results arrived at, may be extended to cover more than the case of two strata, and will help throw more light on the way design effects change as factors like strata means, strata variances, sample sizes and strata proportions vary. Sample designers may -at the design stage- find some guidelines in these findings that can help them determining the expected precision of their intended design.

REFERENCES:

- 1- Kish, L. (1965). *Survey Sampling*. Wiley. New York.
- 2- Kish, L. (1982). Design Effect. *Journal of Encyclopedia of Statistical Sciences*, vol. 2, pp. 347- 348.
- 3- Kish, L. (1995). Methods for design effects. *Journal of Official Statistics*, vol. 11, pp. 55-77.
- 4- Kalton, G. (2005), Michael Brick J., and Thanh Le[^]. Estimating Components of design effects for use in sample design. *Journal of Household Sample Surveys in Developing and Transition Countries*, 96: 95-121.

5- Kalton, G. (1977). Practical methods for estimating survey sampling errors. *Bulletin of the International Statistical Institute*, **47**: 495-514.

6- Curtin L, Carroll M, Dohramann S, Winters F, Lacher D (2002). Extreme Design Effects, Why they occur and what to do about them. *Proceedings of The Joint Statistical Meetings*, New York.

7- David A. Lacher, Lester R. Curtin, Jeffery P. Hughes (2004). Why large Design Effects Can Occur In Complex Sample Designs: Examples from the Nhanes 1999-2000 Survey. *National Center for Health Statistics*. Hyattsville.

8- Cochran, W.G. (1977). *Sampling Techniques*, 3rd Ed. Wiley. New York.

9- Matthew J. Salganik (2006). Variance Estimation, Design Effects, and Sample S for Respondent-Driven Sampling. *Journal of Urban Health*, **83**: 98-112.

10- Douglas C. Montgomery (1997). *Design & Analysis of Experiments*, 5th Ed. Wiley, New York.

11- Klaus Hinkelmann, Oscar Kempthorne (2007). *Design & Analysis of Experiments*, Volume 1. Wiley. New Jersey.

