

**DRAFT****Annex 27****DRAFT BEST PRACTICES EXAMPLES
FOCUSING ON SAMPLE SIZE AND RELIABILITY CALCULATIONS****(Version 01.0)****I Introduction**

1. The clean development mechanism (CDM) Executive Board (hereinafter referred to as the Board) at its fiftieth meeting approved the “General Guidelines for Sampling and Surveys for Small-Scale CDM Project Activities (sampling guideline)”. Further, the Board at its sixtieth meeting agreed to set up a joint task force comprising members of the Methodologies Panel and the Small-Scale Working Group to further work on the issue to develop one set of common sampling guidelines and best practices examples covering large and small-scale projects and programme of activities (PoAs). It further agreed that the scope of the guidelines shall include guidance to designated operational entities (DOEs) on how to review sampling and survey designs in project design documents (PDDs) as well as how to apply sampling to validation/verification work.
2. The Board at its sixty-fifth meeting approved the “Standard for sampling and surveys for CDM project activities and PoAs”.
3. This draft document provides several best practice examples covering large and small-scale project activities. It also provides examples for checking reliability of data collected through sample surveys.
4. This document does not cover the following items, which will be covered in a future document that will be recommended during the year:
 - (a) Methods, if any, to deal with missed reliability targets without compromising conservative estimates for emission reduction;
 - (b) Best practice examples for DOE validation/verification for sampling and surveys.

B. Introductory notes on sample size calculations

5. There are different equations to calculate a required sample size for different situations. Most of the examples in this document are for finite populations such as cook stoves or compact fluorescent lamps (CFLs), but there is also one example for a wastewater treatment plant where the measurements are done for the continuous flow of wastewater.
6. Which equation to use depends on the following:
 - (a) Parameter of interest, for example:
 - (i) A percentage, such as the proportion of annually operating cook stoves;
 - (ii) A numeric value, such as the mean value of operating hours of CFLs, the mean value of dry compressive strength (to check whether the manufactured bricks are of a certain quality);
7. There are other parameters, e.g. ratios, but this guide only covers proportions and means:

**DRAFT**

- (a) Sampling scheme. This document contains the equations and examples using the following five sampling schemes:
 - (i) Simple random sampling;
 - (ii) Systematic sampling;
 - (iii) Stratified random sampling;
 - (iv) Cluster sampling;
 - (v) Multi-stage sampling.
8. There are a number of factors that affect the sample size required, and these are described below:
- (a) The value that the parameter is expected to take, for example:
 - (i) Sampling to see whether 80% of installed cook stoves are still in operation will give a different sample size required than 65%. The same is true for mean values;
 - (b) The amount of variation affects the sample size required. The larger the variation associated with the parameter of interest the larger the sample size required for the same level of confidence and precision;
 - (c) The level of precision (e.g. $\pm 10\%$ of relative value of the parameter's true value) and confidence (e.g. 90% or 95%) in that precision which is desired for determining the parameter also determines the sample size. The higher the required confidence and the narrower the precision the more samples are required.
9. Estimates of the parameter of interest (proportion, mean and standard deviation) are required for sample size calculations. There are different ways to obtain these:
- (a) We may refer to the result of previous studies and use these results;
 - (b) In a situation where we do not have any information from previous studies, we could take a preliminary sample as a pilot and use that sample to provide our estimates;
 - (c) We could use "best guesses" based on the researcher's own experiences.
10. Note that if the standard deviation is unknown but the range (maximum – minimum) is known then a rough "rule of thumb" is that the standard deviation can be estimated as the range divided by 4.
11. Also, for different sampling schemes additional information is required, such as strata estimates rather than just the population information.
12. There are three additional points to make in relation to these sample size calculations:
- (a) If sample size calculations are being performed manually, it is important to retain as many decimal places as relevant, until the final calculated figure is reached. Only then should rounding be carried out. In this document, however, for clarity of presentation the detailed calculations are shown with only a small number of decimal places, although the actual calculations themselves used more than is shown;

**DRAFT**

- (b) Researchers are encouraged to carry out more than one sample size calculation. It is highly unlikely that accurate estimates of the parameters will be available, and so the calculation should be performed for a range of possible estimates (e.g. proportion, or mean and standard deviation), and the largest sample size chosen. This should help to ensure that the sample selected will meet the required reliability criteria;
- (c) The pilot studies that are included in the examples here are deliberately small so that calculations can be illustrated fairly easily. In real-life situations they should be larger than those used here.

1. Sample size calculations - Small-scale examples

13. For all of the small-scale examples below, we require 90% confidence that the margin of error in our estimate is not more than $\pm 10\%$ in relative terms.

Proportional parameter of interest (Cook Stove Project)

14. This section covers sample size calculations based on a proportion (or percentage) of interest being the objective of the project, under four different sampling schemes. Regardless of the sampling scheme used, the following have to be pre-determined in order to estimate the sample size:

- (a) The value that the proportion is expected to take;
- (b) The level of precision, and confidence in that precision (90/10 for all small-scale examples).

15. For all of the cook stove examples below, the proportion of interest is the number of project cook stoves that are still in operation at the end of the third year after the stoves were distributed; it is thought that this proportion is 0.5 (50%). The cook stoves were distributed to 640,000 households, and it has been assumed that 1 household = 1 cook stove.

Example 1 – Simple Random Sampling

16. Suppose that the population is homogenous with respect to the continued use of the cook stoves. Then simple random sampling would be an appropriate method to estimate the proportion of cook stoves still in operation.

The equation to give us the required sample size is:

$$n \geq \frac{1.645^2 NV}{(N-1) \times 0.1^2 + 1.645^2 V} \quad (1)$$

Where:

$$V = \frac{p(1-p)}{p^2}$$

n	Sample size
N	Total number of households (640,000)
p	Our expected proportion (0.50)

DRAFT

1.645 Represents the 90% confidence required

0.1 Represents the 10% relative precision ($0.1 \times 0.5 = 0.05 = 5\%$ points either side of p)

17. Substituting in our values gives:

$$V = \frac{0.5 \times (1 - 0.5)}{0.5^2} = 1 \quad (2)$$

$$n \geq \frac{1.645^2 \times 640,000 \times 1}{(640,000 - 1) \times 0.1^2 + 1.645^2 \times 1} = 270.4 \quad (3)$$

18. Therefore the required sample size is at least 271 households. This assumes that 50% of the cook stoves would be operating. If we changed our prior belief of the underlying true percentage of working stoves p , this sample size would need recalculating.

19. Note that the figure of 271 households means 271 households with data for analysis. If we expected the response rate from the sampled households to be only 80% then we would need to scale up this number accordingly. Thus we would decide to sample $271/0.8 = 339$ households.

20. If we did not scale up our sample size and experienced a response rate of 80% then we would only have 216 ($271 \times 0.8 = 216$) households/cookers with data, and consequently the level of precision would be detrimentally affected. We can calculate the actual level of precision by substituting $n = 216$ into the equation:

$$\frac{1.645^2 \times 640,000 \times 1}{(640,000 - 1) \times \text{precision}^2 + 1.645^2 \times 1} = 216 \quad (4)$$

21. This gives a relative precision of 0.1119, or 11.2% – not the 10% required. So by not adjusting our estimated sample size to take into account the expected response rate we have an increased margin of error.

22. One solution to this could be to take an additional sample of households. This additional sample would need to recruit 69 households which would then, again assuming a response rate of 80%, provide data on 55 households (80% of the 69). Adding these to the existing 216 gives us data for 271 households, the number required to achieve 90/10 reliability.

Approximate equation

23. The equation used above is the exact equation derived from simple random sampling theory. When population sizes are large (or infinite), then an approximate equation can be used, which ignores the actual size of the population (N). The approximate equation for the 90/10 confidence/precision guideline is:

Approximate Equation	Sample size for the above example
----------------------	-----------------------------------

DRAFT

$$\text{Proportion data} \quad n = \frac{1.645^2 V}{0.1^2} \quad \text{Where: } V = \frac{p(1-p)}{p^2} \quad 271 \quad \left(= \frac{1.645^2 \times 1}{0.1^2} \right)$$

Notes on approximate equations

24. As the sample size in this example is large, there is no difference between the sample sizes derived from the exact and approximate equations. However, for smaller populations ($N < 5000$) and small p 's (less than 0.5) there will be a difference.

25. Since the exact equation can be easily calculated, it is recommended that the exact equation be used in preference to the approximate one. It avoids having to decide whether the population size is large enough for it to be possible to use the approximate equation.

26. The scaling-up of the sample size due to non-response will also apply to the approximate equation.

Example 2 – Stratified Random Sampling

27. This time we know that stoves were distributed in four different districts and that the cook stoves are more likely to be still in operation in certain districts compared to others.¹ In this situation we want to take our knowledge about the district differences into account when we do the sampling, and sample separately from each district. Estimates of the proportion of cook stoves still in operation in each district, as well as the population size of each district are required.

District	Number of households with cook stove in district* (g)	Proportion of cook stoves still in operation in district (p)
A	76,021	0.20
B	286,541	0.46
C	103,668	0.57
D	173,770	0.33

* Note that the districts cover all of the population (sum of district populations = total population)

28. The equation for the total sample size is:

$$n \geq \frac{1.645^2 NV}{(N-1) \times 0.1^2 + 1.645^2 V} \quad (5)$$

Where: $V = \frac{SD^2}{\bar{p}^2} = \frac{\text{overall variance}}{\bar{p}^2}$ and \bar{p} is the overall proportion.

29. To then decide on the number of households in the sample that come from each district we could use proportional allocation, where the proportions of units from the different districts in the

¹ If the proportions were expected to be the same in each district then simple random sampling should be used.

DRAFT

sample are the same as the proportions in the population. This gives $n_i = \frac{g_i}{N} \times n$ where $i=1, \dots, k$ and k is the number of districts in the area (in this case 4).

Where:

g_i Size of the i^{th} group (district) where $i=1, \dots, k$

N Population total

30. We use the figures from the table above to calculate the overall variance,² and proportion of cook stoves still in operation.

$$SD^2 = \frac{(g_a \times p_a(1-p_a)) + (g_b \times p_b(1-p_b)) + (g_c \times p_c(1-p_c)) + \dots + (g_k \times p_k(1-p_k))}{N} \quad (6)$$

$$\bar{p} = \frac{(g_a \times p_a) + (g_b \times p_b) + (g_c \times p_c) + \dots + (g_k \times p_k)}{N} \quad (7)$$

Where g_i and N are as above and p_i is the proportion for the i^{th} group (district); $i=1, \dots, k$

Substituting the values from the table into the above equations for SD^2 and \bar{p} gives:

$$SD^2 = \frac{(76021 \times 0.20 \times 0.8) + \dots + (173770 \times 0.33 \times 0.66)}{640000} = 0.23 \quad (8)$$

$$\bar{p} = \frac{(76021 \times 0.20) + (286541 \times 0.46) + (103668 \times 0.57) + (173770 \times 0.33)}{640000} = 0.41 \quad (9)$$

Therefore:

$$V = \frac{SD^2}{\bar{p}^2} = \frac{0.23}{0.41^2} = 1.37 \quad (10)$$

Substituting in V into our sample size equation gives:

$$n \geq \frac{1.645^2 \times 640000 \times 1.37}{(640000 - 1) \times 0.1^2 + 1.645^2 \times 1.37} = 367.0 \quad (11)$$

31. The total sample size required is 367 households. This then needs to be divided up according to the size of each district to get the number of households that should be sampled in each district.

$$\text{General Equation: } n_i = \frac{g_i}{N} \times n \quad (12)$$

² The variance of a proportion is calculated as: $p(1-p)$.

DRAFT

$$\text{District A: } n_a = \frac{76021}{640000} \times 367 = 43.7$$

$$\text{District B: } n_b = \frac{286541}{640000} \times 367 = 164.8$$

$$\text{District C: } n_c = \frac{103668}{640000} \times 367 = 59.6$$

$$\text{District D: } n_d = \frac{173770}{640000} \times 367 = 99.9$$

32. Rounding up the district samples sizes gives the number of households to be sampled in each district, 44 in A, 165 in B, 60 in C, and 100 in D (the sum of these is slightly greater than the total required sample size due to the rounding up of households within each district).

33. Note that these sample sizes do not take into account non-response. If the expected level of response is 75% across all districts then divide each district sample size by 0.75; this will result in larger sample sizes allowing for the non-responders.

Example 3 – Cluster Sampling

34. Now consider a different scenario. The households are not located in different districts. Instead they are ‘clustered’ or grouped into lots of villages. Instead of going to numerous individual households, we want to go to a number of villages and sample every household within each village.

35. For this example the population comprises 120 villages, all of approximately similar size. In order to have some understanding of the proportion of cook stoves still operating and the variation in this proportion between villages, a small preliminary sample has been taken:

Village	Estimated proportion of cook stoves operating in each village
1	0.37
2	0.48
3	0.50
4	0.27
5	0.68
Average (\bar{p})	0.46
Variance SD_B^2	0.024

36. The average (\bar{p}) is just $\frac{0.37 + 0.48 + 0.50 + 0.27 + 0.68}{5} = \frac{2.3}{5} = 0.46$ and the variance between the clusters is:

$$SD_B^2 = \frac{1}{n-1} \sum_{i=1}^{n=5} (p_i - \bar{p})^2 = \frac{(0.37 - 0.46)^2 + (0.48 - 0.46)^2 + \dots + (0.68 - 0.46)^2}{4} = \frac{0.0946}{4} = 0.0237 \quad (13)$$

The equation for the number of villages that need to be sampled is:

$$c \geq \frac{1.645^2 MV}{(M-1) \times 0.1^2 + 1.645^2 V} \quad (14)$$

DRAFT

Where:

$$V = \frac{SD_B^2}{\bar{p}^2} = \frac{\text{variance between clusters (villages)}}{\text{average proportion}}$$

c Number of clusters to be sampled (villages)

M Total number of clusters (villages) – this must encompass the entire population

1.645 Represents the 90% confidence required

0.1 Represents the 10% relative precision required

37. Substituting our values into the above equation gives the number of villages that are required to be sampled as:

$$V = \frac{SD_B^2}{\bar{p}^2} = \frac{0.0237}{0.46^2} = 0.11 \quad (15)$$

$$c \geq \frac{1.645^2 \times 120 \times 0.11}{(120 - 1) \times 0.1^2 + 1.645^2 \times 0.11} = 24.3 \quad (16)$$

38. Therefore we would have to sample every household within 25 randomly selected villages. This approach to sampling assumes that the villages are homogenous. In this example this means that the proportion of cook stoves still operating in a village is independent of any other factors such as district (see example 2 – stratified sampling), economic status, etc. If the proportions are not independent of another factor then cluster sampling within each strata of the factor can be used.

39. Since cluster sampling is dealing with data from whole clusters (villages in this example), non-response at the within-village level (household in this case) is less likely to be an issue, unless there is a high percentage of non-responses within a village. If there are only one or two missing values in a village it is still possible to obtain a usable proportion for that village based on all the other households that did provide data.

Example 4 – Multi-stage Sampling

40. Multi-stage sampling can be thought of as sampling from a number of groups, and then going on to sample units within each group. Continuing with the cook stove example, we want to sample a number of villages and then a number of households within each sampled village.

41. We know that there are 120 villages and there are on average 50 households within each village, of which we plan to sample 10. From a small pilot study we already know the following:

Village	Proportion of cook stoves in operation
A	0.37
B	0.48
C	0.50
D	0.27
E	0.68



DRAFT

42. The equation for the number of villages to be sampled is:

$$c \geq \frac{\frac{SD_B^2}{\bar{p}^2} \times \frac{M}{M-1} + \frac{1}{\bar{u}} \times \frac{SD_w^2}{\bar{p}^2} \times \frac{(\bar{N} - \bar{u})}{(\bar{N} - 1)}}{\frac{0.1^2}{1.645^2} + \frac{1}{M-1} \frac{SD_B^2}{\bar{p}^2}} \quad (17)$$

Where:

c	Number of groups that should be sampled
M	Total number of groups in the population (120 villages)
\bar{u}	Number of units to be sampled within each group (pre-specified as 10 households)
\bar{N}	Average units per group (50 households per village)
SD_B^2	Unit variance (variance between villages)
SD_w^2	Average of the group variances (average within village variation)
\bar{p}	Overall proportion
1.645	Represents the 90% confidence required
0.1	Represents the 10% relative precision

43. Using our table of pilot information we can calculate the unknown quantities for the equation above:

Village	Proportion of cook stoves in operation (p_i)	Variance within village ($p_i(1-p_i)$)
A	0.37	0.2331
B	0.48	0.2496
C	0.50	0.2500
D	0.27	0.1971
E	0.68	0.2176
Average	$\bar{p} = 0.46$	$SD_w^2 = 0.2295$
Variance	$SD_B^2 = 0.0237$	

Where:

\bar{p} is the average proportion of cook stoves, i.e. $\frac{0.37 + \dots + 0.68}{5} = 0.46$

SD_w^2 is the average variance within the villages, i.e. $SD_w^2 = \frac{0.2331 + \dots + 0.2176}{5} = 0.2295$

DRAFT

SD_B^2 is the variance between the village proportions, i.e. the variance between 0.37, 0.48 etc. This can be calculated in the usual way for calculating a variance i.e. using the equation

$$SD_B^2 = \frac{\sum_{i=1}^n (p_i - \bar{p})^2}{n-1} \text{ which gives } SD_B^2 = 0.0237$$

44. Substituting our values into the group sample size equation gives:

$$c \geq \frac{\frac{0.0237}{0.46^2} \times \frac{120}{(120-1)} + \frac{1}{10} \times \frac{0.2295}{0.46^2} \times \frac{(50-10)}{(50-1)}}{\frac{0.1^2}{1.645^2} + \left(\frac{1}{(120-1)} \times \frac{0.0237}{0.46^2} \right)} = 43.4 \quad (18)$$

45. Therefore if we were to sample 10 households from each village we should sample 44 villages for the required confidence/precision.

46. It is usually useful to have this calculation automated so that a series of different u values (the number of units to be sampled in each group) can be used and the effect that this has on the number of groups to be sampled can be observed.

Number of households sampled in each village u	Required number of villages c
5	68
10	44
15	36
20	32
30	28
50	25

47. In this example, by doubling the number of households within each village to be sampled from 10 to 20, we reduce the number of villages that need to be visited from 44 to 32.

48. Note that when u = the average number of households in a village (50), the required sample size is the same as that from cluster sampling as everyone within each village would be sampled. When u is smaller than the average number of households, the number of villages that need to be sampled under multi-stage sampling is greater than that from cluster sampling as not everyone within each village is being sampled.

Mean value parameter of interest (CFL Project)

49. This section covers sample size calculations where the objective of the project relates to a mean value of interest, under four different sampling schemes. For the sample size calculations, regardless of the sampling scheme, we need to know:

- (a) The expected mean (the desired reliability is expressed in relative terms to the mean);
- (b) The standard deviation;

DRAFT

- (c) The level of precision, and confidence in that precision (90/10 for all small-scale examples).

50. The examples below are based on the parameter of interest being average daily CFL usage, which is thought to be 3.5 hours, with a standard deviation of 2.5 hours. The population consists of 420,000 households to which CFLs were distributed; we are assuming that 1 household = 1 CFL.

Example 5 – Simple Random Sampling

51. For simple random sampling to be appropriate we are assuming that CFL usage is homogenous amongst the households.

52. The following equation can be used to calculate the sample size:

$$n \geq \frac{1.645^2 NV}{(N-1) \times 0.1^2 + 1.645^2 V} \quad (19)$$

Where:

$$V = \left(\frac{SD}{mean} \right)^2$$

n	Sample size
N	Total number of households
$mean$	Our expected mean (3.5 hours)
SD	Our expected standard deviation (2.5 hours)
1.645	Represents the 90% confidence required
0.1	Represents the 10% relative precision

$$V = \left(\frac{2.5}{3.5} \right)^2 = 0.51 \quad (20)$$

$$n = \frac{1.645^2 \times 420,000 \times 0.51}{(420,000 - 1) \times 0.1^2 + 1.645^2 \times 0.51} = 138.0 \quad (21)$$

53. Therefore the required sample size is at least 138 households.

54. Note that if we expected the response rate from the sampled households to be only 70% then we would need to scale up the number obtained above accordingly. Thus we would decide to sample $138/0.7 = 198$ households.

Approximate equation

55. The equation used above is the exact equation. When population sizes are large (or infinite), then approximate equations can be used, which ignore the actual size of the population (N).

56. The approximate equation follows the 90/10 confidence/precision guideline:

DRAFT

	Approximate Equation	Sample size for the above example
Mean value data	$n = \frac{1.645^2 V}{0.1^2}$ Where: $V = \left(\frac{SD}{mean} \right)^2$	$138 \left(= \frac{1.645^2 \times 0.51}{0.1^2} \right)$

57. Please see the ‘Approximate equation’ section under **i) Cook Stove Project - Proportional parameter of interest, Example 1: Simple Random Sampling** for notes relating to approximate equations.

Example 6 – Stratified Random Sampling

58. The key to this example is that, unlike under simple random sampling, it is not assumed that the population is homogeneous – different parts of the population are expected to have different CFL usage averages.

59. Suppose that the CFLs were distributed in different districts in which each has a different CFL usage pattern (due to district economic backgrounds). We are now interested in sampling users of CFLs from all the districts to ensure all areas are well represented.

60. Each district has the following number of households, and mean and standard deviation CFL usage:

District	Number of households in district given a CFL	Mean (hours)	Standard deviation (hours)
A	146,050	3.2	1.9
B	104,474	2.4	0.8
C	38,239	4.5	1.6
D	74,248	1.6	1.7
E	56,989	2.3	0.7

61. The total sample size of households across all five districts is:

$$n \geq \frac{1.645^2 \times NV}{(N-1) \times 0.1^2 + 1.645^2 V} \quad (22)$$

Where:

$$V = \left(\frac{SD}{mean} \right)^2$$

SD Is the overall standard deviation, and

mean Is the overall mean.

62. Using the data in the table above we can estimate the overall mean and standard deviation. Both equations are weighted according to the total number of households in each district.

63. Overall Standard Deviation:

DRAFT

$$SD = \sqrt{\frac{(g_a \times SD_a^2) + (g_b \times SD_b^2) + (g_c \times SD_c^2) + \dots + (g_k \times SD_k^2)}{N}} \quad (23)$$

Where:

SD Weighted overall standard deviation
 SD_i Standard deviation of the i^{th} group where $i=1, \dots, k$, (note that these are all squared – so the group size is actually being multiplied by the group variance)

g_i Size of the i^{th} group where $i=1, \dots, k$

N Population total

$$mean = \frac{(g_a \times m_a) + (g_b \times m_b) + (g_c \times m_c) + \dots + (g_k \times m_k)}{N} \quad (24)$$

Where:

$mean$ Weighted overall mean

m_i Mean of the i^{th} group where $i=1, \dots, k$

64. Substituting the values from our example into the above expressions gives:

$$SD = \sqrt{\frac{(146050 \times 1.9^2) + (104474 \times 0.8^2) + \dots + (56989 \times 0.7^2)}{420000}} = 1.49 \quad (25)$$

$$mean = \frac{(146050 \times 3.2) + \dots + (56989 \times 2.3)}{420000} = 2.71 \quad (26)$$

65. Substituting these values into the equation for V gives:

$$V = \left(\frac{SD}{mean} \right)^2 = \left(\frac{1.49}{2.71} \right)^2 = 0.3 \quad (27)$$

And hence, for the sample size:

$$n = \frac{1.645^2 \times 420,000 \times 0.3}{(420,000 - 1) \times 0.1^2 + 1.645^2 \times 0.3} = 81.7 \quad (28)$$

66. This example assumes proportional allocation, which means that the number of households we want to sample from each district is proportional to the size of the district within the population.

The equation for each district sample size is: $n_i = \frac{g_i}{N} \times n$ (29)

$$\text{District A: } n_a = \frac{146050}{420000} \times 82 = 29$$

$$\text{District B: } n_b = \frac{104474}{420000} \times 82 = 21$$

DRAFT

$$\text{District C: } n_c = \frac{38239}{420000} \times 82 = 8$$

$$\text{District D: } n_d = \frac{74248}{420000} \times 82 = 15$$

$$\text{District E: } n_e = \frac{56989}{420000} \times 82 = 12$$

67. The summation of these district sample sizes (29+21+8+15+12=85) is slightly greater than that calculated from the total sample size equation (82) above due to rounding.

68. As with previous examples, the sample sizes above need to be scaled up to take into account any non-response expected.

Example 7 – Cluster Sampling

69. Suppose CFLs were distributed to households in 50 villages. Instead of sampling from the whole population of households with CFLs, we sample a number of villages (villages=clusters), and then collect data from all households within the villages.

70. The equation used to give us the required number of clusters, c , to sample is:

$$c \geq \frac{1.645^2 MV}{(M-1) \times 0.1^2 + 1.645^2 V} \quad (30)$$

Where:

$$V = \left(\frac{SD}{\text{Cluster mean}} \right)^2$$

M Total number of clusters (50 villages)

1.645 Represents the 90% confidence required

0.1 Required precision

71. To perform the calculations we need information about CFL usage at the village level, rather than at the household level. If such information does not already exist, we could possibly collect it in a pilot study. The example here assumes that data are available from a pilot study on five villages.

Village	Total usage across all households in the village ³
A	30458
B	27667
C	31500
D	28350
E	19125

³ In the pilot study these totals may be derived from collecting data on all households in the village, or else by taking a sample of households in the village and scaling up from the sample to all households.

DRAFT

72. Calculating the mean and standard deviation of these figures gives:

$$\text{Cluster mean } (\bar{y}) = \frac{1}{n} \sum_{i=1}^n y_i = \frac{30458 + 27667 + \dots + 19125}{5} = 27420 \quad (31)$$

$$SD = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2} \text{ where the } y_i \text{ are the total usages for the villages.}$$

$$SD_B^2 = 23902660 \text{ and so } SD_B = 4889$$

73. These statistics (i.e. mean and SD of data) are easily produced using statistical software.

74. Substituting these values into the equation gives the required number of clusters, i.e. villages as:

$$V = \left(\frac{4889}{27420} \right)^2 = 0.03 \quad (32)$$

$$c \geq \frac{1.645^2 \times 50 \times 0.03}{(50-1) \times 0.1^2 + 1.645^2 \times 0.03} = 7.5 \quad (33)$$

75. So we need to sample eight villages to satisfy the 90/10 confidence/precision criterion. Once a village is selected, all households in the selected village should be sampled.

76. The above equation assumes that CFL usage in a village is independent of any other factors, such as economic status. If CFL usage was expected to vary according to another factor then cluster sampling can be used within each level of the factor.

Example 8 – Multi-stage Sampling

77. Multi-stage sampling combines the cluster and simple random sampling approaches in a two-stage sampling scheme which enables us to randomly select some groups (villages) and then randomly sample some units (households) within those groups (villages). As with simple random sampling and cluster sampling, we are assuming homogeneity across villages in the usage of CFLs. We know that the 420,000 households are in 50 villages.

78. Let us start by assuming that we want to sample 10 households in each village. In general terms we will call this number u (for units).

79. In order to perform a sample size calculation we need information on:

- (i) The variation between households within the villages;
- (ii) The variation between villages;
- (iii) The average household usage;
- (iv) The average usage at the village level.

80. A previous study had provided data for households in five villages, and the results are summarized below. Note that not all villages in this example are exactly the same size.

DRAFT

CFL average daily usage (hours)				
Village	Number of households	Mean usage ⁴ per household in a village	Total usage across all households	Standard deviation ⁵ (between households within villages)
A	8500	3.58	30458	2.60
B	8300	3.33	27667	2.70
C	8400	3.75	31500	0.66
D	8100	3.50	28350	0.75
E	8500	2.25	19125	1.50
Total number of households	41800			
Overall mean usage per household		3.28		
Mean usage per village			27420	
SD_B = Standard deviation between villages (SD of the total usage column)			4889	
SD_W = Average within village standard deviation				1.86

81. In the table above, the overall mean CFL usage is the average usage for a household, i.e.

$$\text{Overall mean} = \frac{30458 + 27667 + \dots + 19125}{41800} = 3.28$$

82. The cluster or village mean CFL usage is the average usage for village, i.e.

$$\text{Cluster mean} = \frac{30458 + 27667 + \dots + 19125}{5} = 27420$$

83. SD_W^2 is the average of the variances between households within the villages. Its square root (SD_W) is the average within village standard deviation. The equation for SD_W^2 is:

$$SD_W^2 = \frac{8500 \times 2.60^2 + \dots + 8500 \times 1.50^2}{41800} = 3.48 \text{ and so } SD_W = 1.86$$

SD_B^2 is the variance between the village total usages and its square root is the standard deviation between villages. It can be calculated using the usual equation for a variance, i.e.

$$SD_B^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1} \text{ where the } y_i \text{ are the total usages for the villages.}$$

$$SD_B^2 = 23902660 \text{ and so } SD_B = 4889$$

⁴ This can be a mean from all households or a mean from a sample of households.

⁵ And this can be a standard deviation based on all households or a sample of households.

DRAFT

84. We have pre-specified that we want to sample 10 households within each village, so we need to calculate how many villages need to be sampled given a 90/10 confidence/precision criterion is required:

$$c \geq \frac{\left(\frac{SD_B}{Clustermean} \right)^2 \times \left(\frac{M}{M-1} \right) + \left(\frac{1}{u} \right) \times \left(\frac{SD_W}{Overallmean} \right)^2 \left(\frac{\bar{N}-u}{\bar{N}-1} \right)}{\left(\frac{0.1}{1.645} \right)^2 + \frac{1}{M-1} \left(\frac{SD_B}{Clustermean} \right)^2} \quad (34)$$

Where:

M Total number of groups (50 villages)

\bar{N} Average number of units per group (approximately 8,400 households per village)

u Number of units that have been pre-specified to be sampled per group (pre-specified number of households to be sampled in each village = 10)

1.645 Represents the 90% confidence required

0.1 Required precision

$$c \geq \frac{\left(\frac{4889}{27420} \right)^2 \times \left(\frac{50}{50-1} \right) + \left(\frac{1}{10} \right) \times \left(\frac{1.86}{3.28} \right)^2 \left(\frac{8400-10}{8400-1} \right)}{\left(\frac{0.1}{1.645} \right)^2 + \frac{1}{50-1} \left(\frac{4889}{27420} \right)^2} = 14.9 \quad (35)$$

85. Therefore if we were to sample 10 households from each village we should sample 15 villages for the required confidence/precision.

86. It is usually useful to have this calculation automated so that a series of different u values (the number of units to be sampled in each group) can be used and the effect that this has on the number of groups to be sampled can be observed.

Number of households sampled in each village u	Required number of villages c
5	23
10	15
15	13
20	12
25	11

87. Compared to the cluster sampling example, more villages are required for the multi-stage sampling scheme because fewer households are being sampled within each village.

88. Note that in the above example the villages are of slightly different sizes. In practice this is likely to be the case, although the actual sizes may not always be known. This is not critical to the sample size calculation. What is important is that sensible estimates of the mean and standard

**DRAFT**

deviation at both the cluster level (village level) and unit level (household level) are used in the calculation.

Mean value parameter of interest (Brick Project)

89. This section covers an example sample size calculation based on systematic sampling where the objective of the project relates to a mean value of interest.

As for all mean value parameters of interest examples we need to know:

- (a) The expected mean (the desired reliability is expressed in relative terms to the mean);
- (b) The standard deviation;
- (c) The level of precision, and confidence in that precision (90/10 for all small-scale examples).

90. The following example is based on assessing whether bricks are of a minimum quality after manufacture; dry compressive strength has been identified as a suitable measurement of quality. Prior information gives us a mean dry compressive strength of 158kg/cm² with a standard deviation of 65kg/cm².

Example 9 – Systematic Sampling

91. This example is based on a manufacturing process; we want to systematically sample every nth brick from the production line of 500,000 bricks per year. We wish to know how many bricks should be sampled to ensure an average dry compressive strength of 158kg/cm², with 90/10 confidence/precision.

92. The sample size equation for a required 90/10 confidence/precision is:

$$n \geq \frac{1.645^2 V}{0.1^2} \quad (36)$$

Where:

$$V = \left(\frac{SD}{mean} \right)^2$$

93. Substituting in mean and standard deviation from above gives:

$$V = \left(\frac{65}{158} \right)^2 = 0.17 \quad (37)$$

$$n \geq \frac{1.645^2 \times 0.17}{0.1^2} = 45.8 \quad (38)$$

94. Therefore, we should take 46 samples to gain the required levels of confidence and precision. Given that 500,000 bricks are manufactured each year and we want to take 46 samples,

**DRAFT**

we should sample 1 brick for every N/n bricks produced – that is 1 brick for every 10,917 (= 500,000 / 46) bricks.

95. To make sure our sample is random, we randomly choose a starting point (starting brick) between 1 and 10,917 and use this brick as our first sample – for example brick 6505. We continue sampling by taking every 10,917th brick, so our second sample would be brick 6505 + 10,917 = 17,422, the third brick sampled would be 17,422 + 10,917 = 28,339, etc. For the sake of practicality it might be easier to sample every 10,000th brick; instead of every 10,917th, this would give a slightly larger sample size.

*Measurements in biogas projects*Example 10

96. A survey will be carried out to estimate the mean chemical oxygen demand (COD) at a wastewater plant. The wastewater is a continuous flow of water that leaves the plant. A 500 ml sample of water will be extracted (from plant inlet) from the continuous flow of wastewater on a regular basis throughout the year and a single measurement of COD (mg/L) made on each sample.

97. This form of sampling, i.e. on a regular basis, possibly with a random start date, is systematic sampling.

98. The wastewater system has been in place for some time, and is considered to be stable in terms of the way it is functioning. The COD for the inlet is thought to be at a constant level throughout the year (apart from random variation)⁶.

99. Previous work where measurements were taken on a regular basis suggested that the mean COD for untreated water is likely to be about 31,750 mg/L and the standard deviation (*SD*) in the order of 6,200 mg/L.

100. Since the wastewater is flowing continuously, the study population can be thought of as all possible 500 ml water samples in a whole year – so large as to be almost infinite. The sample size calculation no longer needs inclusion of the finite population size (i.e. *N*).

101. If the sampling times are sufficiently far apart the data can be regarded as a set of independent observations and treated as a simple random sample. The number of COD measurements that are required to meet the 90/10 reliability is:

$$n = \left(\frac{t_{n-1} \times SD}{0.1 \times \text{mean}} \right)^2 \quad (39)$$

102. Where t_{n-1} is the value of the t-distribution for 90% confidence when the sample size is n .⁷ However, the sample size is not yet known, and so a first step is to use the value for 90% confidence when the sample is large, i.e. 1.645, and then refine the calculation.

⁶ In reality, temporal fluctuations (daily, weekly, seasonally, etc.) both in the wastewater flow and COD concentration should be taken into account when taking samples.

⁷ This is indicated by the subscript ($n-1$) which is called the degrees of freedom for the t-value.

DRAFT

$$n = \left(\frac{1.645 \times SD}{0.1 \times \text{mean}} \right)^2 \quad (40)$$

103. This gives $n = \left(\frac{1.645 \times 6200}{0.1 \times 31750} \right)^2 = 10.3$ which rounds up to 11.

104. The calculation now needs to be repeated using the t-value for 90% confidence and n=11.

105. The exact figure for this t-value can be acquired from any set of general statistical tables or using standard statistical software. For a sample size of 11 the value is 1.812.

106. The calculation now gives $n = \left(\frac{1.812 \times 6200}{0.1 \times 31750} \right)^2 = 12.5$ which rounds up to 13.

107. The process should be iterated until there is no change to the value of n. Here the repeat calculation would have a t-value of 1.782 and the calculation would yield n = 12.11, which would be rounded up to 13. The sample size calculation suggests that sampling every four weeks should be sufficient for 90/10 reliability.

Other calculated sample sizes

108. The above is a relatively simple example, and not all situations will yield values as neat as “once a month” or “once every four weeks”. For instance:

- Had the calculation indicated that 48 measurements should be taken, one would most likely decide to sample weekly for the whole year;
- If the calculation had indicated 16 samples were required, then one might decide to sample every three weeks. Alternatively, since this may not be an easy schedule to comply with, one might choose to sample every two weeks. This now gives us a total of 26 samples which should ensure that the data, when collected and analysed, have more than adequate precision (assuming, of course, that the figures for the mean and the standard deviation that were used in the sample size calculation were good reflections of the true situation).

109. Instead of trying to follow “unworkable” schedules, it may be more sensible to use the following simplifications:

No. of measurements determined from sample size calculation	Proposed schedule
Less than or equal to 12	Monthly
13–17	Every three weeks
18–26	Every two weeks
26–51	Every week
More than 52	Twice a week

Understanding variation

DRAFT

110. The above example illustrates the sample size calculation using an absolute figure for the standard deviation. However, sometimes researchers have difficulty providing a figure for the standard deviation, but they can express it in relative terms. For instance, when asked about the variation in COD in this wastewater example, the researcher may describe it as 20%.

111. The coefficient of variation (CV) is a summary measure which describes variability in terms of the mean. The actual equation is $CV = \frac{SD}{mean}$. It is sometimes multiplied by 100, in which case it is describing the standard deviation as a percentage of the mean. The sample size equation on the previous page can now be written as $n = \left(\frac{t_{n-1} \times CV}{0.1} \right)^2$ where t_{n-1} is the value of the t-distribution for 90% confidence for a sample of n measurements. Again the value of 1.645 would be used instead of a t-value for the first step in the calculation; and so in this example the first step would be $n = \left(\frac{1.645 \times 0.2}{0.1} \right)^2 = 10.8$.

Sample size calculations - Large-scale examples

112. For the large-scale examples we require 95% confidence that the margin of error in our estimate is not more than $\pm 10\%$ in relative terms.

Proportional parameter of interest (Transport Project)

113. This section covers sample size calculations based on a proportion (or percentage) of interest being the objective of the project, under four different sampling schemes. Regardless of the sampling scheme used, the following have to be pre-determined in order to estimate the sample size based on a proportion:

- (a) The value that the proportion is expected to take;
- (b) The level of precision, and confidence in that precision (95/10 for all large-scale examples).

114. The examples relate to passengers travelling on the transport project in Bogota. It is known that 1,498,630 passengers use the project for transportation every day; the parameter of interest is the proportion of these passengers that would previously have travelled by bus, thought to be 45%.

Example 11 – Simple Random Sampling

115. Assuming that the proportion of interest is homogenous. The equation for the sample size required under simple random sampling is:

$$n \geq \frac{1.96^2 NV}{(N-1) \times 0.1^2 + 1.96^2 V} \quad (41)$$

Where:

$$V = \frac{p(1-p)}{p^2}$$



DRAFT

n	Sample size with finite population correction
N	Total number of passengers per day
p	Our estimated proportion (45%)
1.96	Represents the 95% confidence required
0.1	Required precision

116. Substituting our values into the above equation we get:

$$V = \frac{0.45 \times (1 - 0.45)}{0.45^2} = 1.22 \quad (42)$$

$$n \geq \frac{1.96^2 \times 1,498,630 \times 1.22}{(1,498,630 - 1) \times 0.1^2 + 1.96^2 \times 1.22} = 469.4 \quad (43)$$

117. The required sample size is at least 470 passengers to get an estimated proportion of passengers that would have previously travelled by bus with 95/10 confidence/precision. Note that the sample size will change depending on the estimated proportion value.

118. The above sample size does not take into account any non-responders, that is passengers who do not respond to the question. If we expect that 90% of passengers will respond (and 10% will not) then we should increase the sample size by dividing the sample size calculated above by the expected level of response: $470/0.9 = 523$. To account for a non-response of 10% we should sample 523 passengers.

119. Please see the ‘Approximate equation’ section under **i) Cook Stove Project - Proportional parameter of interest, Example 1: Simple Random Sampling** for notes relating to approximate equations. Note that 1.96 should be used in place of 1.645 to account for the increased confidence required for the large-scale projects.

Example 12 – Stratified Random Sampling

120. Suppose that we believe the proportion of transport project passengers that would have previously used the bus would vary between the eight different zones in which the project operates. We would like to make sure that when we do our sampling, our sample includes a representative proportion of passengers from each zone. To calculate the sample size, estimates of the number of passengers and proportion that would have previously travelled by bus within each zone are required.

Zone	Number of passengers within each zone (g)	Estimated proportion of passengers that would have used a bus (p)
A	19865	0.43
B	21358	0.57
C	301245	0.4
D	65324	0.71

DRAFT

E	654832	0.32
F	50213	0.46
G	12489	0.26
H	373304	0.68

121. The equation for the total sample size is:

$$n \geq \frac{1.96^2 NV}{(N-1) \times 0.1^2 + 1.96^2 V} \quad (44)$$

$$\text{Where: } V = \frac{SD^2}{\bar{p}^2} = \frac{\text{overall variance}}{(\text{overall proportion})^2}$$

122. To then decide on the number of passengers in the sample that come from each zone we could use proportional allocation, where the proportions of units from the different zones in the sample is the same as the proportions in the population. This gives:

$$n_i = \frac{g_i}{N} \times n \quad \text{where } i = 1, \dots, k \quad \text{where } k \text{ is the number of zones (in this case 8).}$$

Where:

g_i Size of the i^{th} group (district) where $i=1, \dots, k$

N Population total

123. Using the figures from the table we can calculate the overall variance,⁸ and overall proportion:

$$SD^2 = \frac{(g_a \times p_a (1 - p_a)) + (g_b \times p_b (1 - p_b)) + (g_c \times p_c (1 - p_c)) + \dots + (g_k \times p_k (1 - p_k))}{N} \quad (45)$$

$$\bar{p} = \frac{(g_a \times p_a) + (g_b \times p_b) + (g_c \times p_c) + \dots + (g_k \times p_k)}{N} \quad (46)$$

Where g_i and N are as above and p_i is the proportion of the i^{th} group (district) where $i=1, \dots, k$.

$$SD^2 = \frac{(19865 \times 0.43 \times 0.57) + (21358 \times 0.57 \times 0.43) + \dots + (373304 \times 0.28 \times 0.72)}{1,498,630} = 0.22 \quad (47)$$

$$\bar{p} = \frac{(19865 \times 0.43) + (21358 \times 0.57) + \dots + (373304 \times 0.28)}{1,498,630} = 0.45 \quad (48)$$

124. Therefore:

⁸ The variance of a proportion is calculated as: $p(1-p)$.

DRAFT

$$V = \frac{SD^2}{p} = \frac{0.22}{0.45^2} = 1.09 \quad (49)$$

125. Substituting in our V gives:

$$n \geq \frac{1.96^2 \times 1,498,630 \times 1.09}{(1,498,630 - 1) \times 0.1^2 + 1.96^2 \times 1.09} = 419.6 \quad (50)$$

126. The total sample size required is 420 passengers. The next step is to divide this total sample size up according to the size of each zone to get the number of passengers to be sampled within each zone.

$$n_i = \frac{g_i}{N} \times n \quad (51)$$

General Equation:

$$\text{Zone A: } n_a = \frac{19865}{1,498,630} \times 420 = 5.6$$

$$\text{Zone B: } n_b = \frac{21358}{1,498,630} \times 420 = 6.0$$

$$\text{Zone C: } n_c = \frac{301245}{1,498,630} \times 420 = 84.4$$

$$\text{Zone D: } n_d = \frac{65324}{1,498,630} \times 420 = 18.3$$

$$\text{Zone E: } n_e = \frac{654832}{1,498,630} \times 420 = 183.5$$

$$\text{Zone F: } n_f = \frac{50213}{1,498,630} \times 420 = 14.1$$

$$\text{Zone G: } n_g = \frac{12489}{1,498,630} \times 420 = 3.5$$

$$\text{Zone H: } n_h = \frac{373304}{1,498,630} \times 420 = 104.6$$

127. Rounding up the zone sample sizes gives the number of passengers to be sampled in each zone (the sum of these is slightly greater than the required sample size due to the rounding up of passengers within each zone). The sample sizes required vary so much between the zones because the number of passengers in each zone is so different.

128. Note that these sample sizes do not take into account non-response. If the expected level of response is 85% across all zones then divide each zone sample size by 0.85. This will result in larger sample sizes allowing for the non-responders.

Example 13 – Cluster Sampling

129. Instead of sampling individual passengers, it has been decided that buses (clusters) are going to be sampled and then every passenger on each of the selected buses will be asked if they travelled by bus prior to the project. To calculate the sample size we require the number of clusters that make up the population, that is the number of buses that carry the transport project passengers; for this example we will assume 12,000 buses. We also need estimated proportions of passengers that would have travelled by bus prior to the project from a number of buses; for this example we have previously sampled four buses and the proportions were:

DRAFT

Bus	Estimated Proportion
1	0.37
2	0.46
3	0.28
4	0.52
Average (\bar{p})	0.4075
Variance (SD_B^2)	0.011

130. The equation for the number of buses that need to be sampled is:

$$c \geq \frac{1.96^2 MV}{(M-1) \times 0.1^2 + 1.96^2 V} \quad (52)$$

Where:

$$V = \frac{SD_B^2}{\bar{p}^2} = \frac{\text{Variance between clusters(buses)}}{\text{Average proportion over clusters}}$$

131. The average proportion is just $\frac{0.37 + 0.46 + 0.28 + 0.52}{4} = \frac{1.63}{4} = 0.41$ and the variance between the clusters is:

$$SD_B^2 = \frac{1}{n-1} \sum_{i=1}^{n=5} (p_i - \bar{p})^2 = \frac{(0.37 - 0.4065)^2 + (0.46 - 0.4065)^2 + \dots + (0.52 - 0.4065)^2}{3} = 0.0110 \quad (53)$$

Where

- c Number of clusters to be sampled (buses)
- M Total number of clusters (buses) - this must encompass the entire population
- 1.96 Represents the 95% confidence required
- 0.1 Represents the 10% relative precision

132. Substituting our values into the above equation gives:

$$V = \frac{SD_B^2}{\bar{p}^2} = \frac{0.0110}{0.41^2} = 0.07 \quad (54)$$

$$c \geq \frac{1.96^2 \times 12000 \times 0.0664}{(12000 - 1) \times 0.1^2 + 1.96^2 \times 0.0664} = 25.5 \quad (55)$$

133. Therefore we would have to sample every passenger on 26 randomly selected buses.

134. This approach to sampling assumes that the population is homogenous. In this example this means that the proportion of passengers that would have previously travelled by bus is independent of any other factors such as zones (see example 12 – stratified sampling), economic status, etc. If

DRAFT

the proportion of passengers that would have previously travelled by bus is expected to be different for different zones, then cluster sampling should be used within each zone.

135. Non-response is unlikely to be a problem when using cluster sampling, unless the number of individuals within a cluster could be 0 (a bus with no passengers). If it is thought that it could be a problem then the sample size should be scaled up accordingly.

Example 14 – Multi-stage Sampling

136. Instead of sampling every passenger on a number of selected buses, suppose we only want to sample a number of passengers on each bus. This can be thought of as multi-stage sampling as we are sampling a number of buses (groups), and then going on to sample units (passengers) within each group.

137. We know that there are 12,000 buses and there are on average 30 passengers on each bus, of which we plan to sample 15. From a small pilot study we already know the following:

Bus	Proportion of passengers that would have travelled by bus
1	0.37
2	0.46
3	0.28
4	0.52

138. The equation for the number of buses to be sampled is:

$$c \geq \frac{\frac{SD_B^2}{\bar{p}^2} \times \frac{M}{M-1} + \frac{1}{\bar{u}} \times \frac{SD_w^2}{\bar{p}^2} \times \frac{(\bar{N} - \bar{u})}{(\bar{N} - 1)}}{\frac{0.1^2}{1.96^2} + \frac{1}{M-1} \times \frac{SD_B^2}{\bar{p}^2}} \quad (56)$$

Where:

c	Number of groups that should be sampled
M	Total number of groups in the population (12,000 buses)
\bar{u}	Number of units to be sampled within each group (pre-specified as 15 passengers)
\bar{N}	Average units per group (30 passengers on each bus)
SD_B^2	Unit variance (variance between buses)
SD_w^2	Average of the group variances (average within bus variation)
\bar{p}	Overall proportion
1.96	Represents the 95% confidence required
0.1	Represents the 10% absolute precision

139. Using our table of pilot information we can calculate the unknown quantities for the equation above:



DRAFT

Bus	Proportion of passengers that would have used a bus p_i	Variance within bus $p_i(1-p_i)$
1	0.37	0.2331
2	0.46	0.2484
3	0.28	0.2016
4	0.52	0.2496
Variance	$SD_B^2 = 0.0110$	
Average	$\bar{p} = 0.41$	$SD_W^2 = 0.2332$

Where:

\bar{p} is the average proportion of passengers who travel by bus, i.e. $\bar{p} = \frac{0.37 + \dots + 0.52}{4} = 0.41$

SD_W^2 is the average variance between passengers on a bus, i.e.

$$SD_W^2 = \frac{0.2331 + \dots + 0.2496}{4} = 0.2332$$

SD_B^2 is the variance between the bus proportions, i.e. the variance between 0.37, 0.48, etc. This can be calculated in the usual way for calculating a variance, i.e. using the equation

$$SD_B^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1} \text{ which gives } SD_B^2 = 0.0110$$

140. Substituting our values into the group sample size equation gives:

$$c \geq \frac{\frac{0.0110}{0.41^2} \times \frac{12000}{(12000-1)} + \frac{1}{15} \times \frac{0.2332}{0.4075^2} \times \frac{(30-15)}{(30-1)}}{\frac{0.1^2}{1.96^2} + \left(\frac{1}{(12000-1)} \times \frac{0.0110}{0.41^2} \right)} = 44.0 \quad (57)$$

141. Therefore if we were to sample 15 passengers from each bus we should sample 44 buses for the required confidence/precision. The table below gives the number of buses required (c) when choosing to sample different numbers of passengers from each bus (u).

Number of passengers sampled on each bus u	Required number of buses c
5	119
10	63
15	44
20	35
30	26

**DRAFT**

142. The required sample size from the cluster sampling scheme example was 26; this is the same as the sample size required under the multi-stage sampling scheme when $u=30$ (the number of passengers sampled from each bus). This is because we assumed an average of 30 passengers on each bus in the calculations, so when we take u =assumed average passengers the two sampling schemes are the same.

Mean value parameter of interest (Transport Project)

143. This section covers sample size calculations where the objective of the project relates to a mean value of interest, under four different sampling schemes. For the sample size calculations, we need to know:

- (c) The expected mean (the desired reliability is expressed in relative terms to the mean);
- (d) The standard deviation;
- (e) The level of precision, and confidence in that precision (95/10 for all large-scale examples).

144. The parameter of interest in the examples below is average journey length (km) of people who travel by car, whether this is a domestic car or a taxi, and for people who travel by bus.

Example 15 – Simple Random Sampling

145. Suppose we are interested in the average distance (km) of car journeys in Bogota in a day (including both domestic cars and taxis), and we assume that the journeys are homogenous. We know that 2,000,000 journeys are made each day and believe that the mean is 8 km with a standard deviation of 3.5 km. Using simple random sampling the sample size equation is:

$$n = \frac{1.645^2 NV}{(N-1) \times 0.1^2 + 1.645^2 V} \quad (58)$$

Where:

$$V = \left(\frac{SD}{mean} \right)^2$$

n Sample size

N Total number of journeys (2,000,000)

$mean$ Expected mean journey length (8 km)

SD Expected standard deviation for the journey length (3.5 km)

1.96 Represents the 95% confidence required

0.1 Represents the 10% precision

$$V = \left(\frac{3.5}{8} \right)^2 = 0.19 \quad (59)$$

DRAFT

$$n = \frac{1.96^2 \times 2,000,000 \times 0.19}{(2,000,000 - 1) \times 0.1^2 + 1.96^2 \times 0.19} = 73.5 \quad (60)$$

146. Therefore the required sample size is at least 74 journeys to find the average journey length with 95% confidence and a 10% relative margin of error.

147. The calculation above does not take into account non-response. If a level of non-response is expected within the sample then the sample size should be scaled up accordingly. For example if we expected 95% of the people on journeys sampled to respond then we should take this into account and plan to sample $74/0.95 = 78$ journeys instead of 74.

148. There is an approximate equation for this sample size calculation. Please see the ‘Approximate equation’ section under **ii) CFL Project – Mean value parameter of interest, Example 5: Simple Random Sampling** for notes relating to approximate equations. Note that 1.96 should be used in place of 1.645 to account for the increased confidence required in the large-scale projects.

Example 16 – Stratified Random Sampling

149. The fundamental aspect of this sampling scheme is that the average journey length differs between domestic cars and taxis (it is not homogenous as assumed in the previous example). Because we know that the type of vehicle affects the journey distance, we want to make sure that we sample a representative number of domestic cars and taxis. A summary of each stratification group is given below:

Stratification group	Number of journeys per day	Mean (km)	Standard Deviation (km)
Domestic Car	1,595,169	9	3.7
Taxi	982,224	7	2.5

150. Using the data in the table above we can estimate the overall mean and standard deviation:

Overall mean:

$$mean = \frac{(g_a \times m_a) + (g_b \times m_b) + (g_c \times m_c) + \dots + (g_k \times m_k)}{N} \quad (61)$$

Where:

$mean$	Weighted overall mean
g_i	Size of the i^{th} group where $i=1, \dots, k$
m_i	Mean of the i^{th} group where $i=1, \dots, k$
N	Population total

Substituting the values from our example into the above expression gives:

$$mean = \frac{(1595169 \times 9) + (982224 \times 7)}{(1595169 + 982224)} \quad (62)$$

$$mean = 8.2$$

DRAFT

Overall Standard Deviation:

$$SD = \sqrt{\frac{(g_a \times SD_a^2) + (g_b \times SD_b^2) + (g_c \times SD_c^2) + \dots + (g_k \times SD_k^2)}{N}} \quad (63)$$

Where:

SD Weighted overall standard deviation

SD_i Standard deviation of the i^{th} group where $i=1, \dots, k$, (note that these are all squared – so the group size is actually being multiplied by the group variance)

151. Using the values from our example gives:

$$SD = \sqrt{\frac{(1595169 \times 3.7^2) + (982224 \times 2.5^2)}{(1595169 + 982224)}} \quad (64)$$

$$SD = 3.3$$

152. The sample size equation uses the overall mean and standard deviation calculated above:

$$n \geq \frac{1.96^2 \times NV}{(N-1) \times 0.1^2 + 1.96^2 V} \quad (65)$$

153. Substituting in the values from our examples gives:

$$V = \left(\frac{SD}{mean} \right)^2 = \left(\frac{3.3}{8.2} \right)^2 = 0.16 \quad (66)$$

$$n = \frac{1.96^2 \times 2577393 \times 0.16}{(2577393 - 1) \times 0.1^2 + 1.96^2 \times 0.16} = 61.4 \quad (67)$$

154. This gives us the total number of journeys that should be sampled across both vehicle types. The section below assumes proportional allocation – which means that the number of journeys we want to sample from each vehicle type is proportional to the number of journeys made by each vehicle type within the population.

$$n_i = \frac{g_i}{N} \times n$$

General equation: (68)

$$\text{Domestic cars: } n_{Car} = \frac{1592169}{2577393} \times 62 = 38.4 \quad \text{Taxis: } n_{Taxi} = \frac{982224}{2577393} \times 62 = 23.6$$

155. Rounding these figures results in a sample consisting of 39 domestic car journeys, and 24 taxi journeys. The summation of these group sample sizes ($39 + 24 = 63$) is slightly greater than that calculated from the equation above due to rounding.

156. The sample size calculated above assumes 100% response, and therefore needs to be scaled up where non-response is likely to occur.

DRAFT

157. This sample size is smaller than that from simple random sampling in this example. This is due to the standard deviations within strata being smaller than the standard deviation across the whole population (which is usually the case).

Example 17 – Cluster Sampling

158. Now consider a different scenario. The parameter of interest is the average journey length of people who take local buses. Instead of sampling numerous individual passengers we would like to sample everyone from a few buses (clusters). Knowing that there are 12,000 local buses in Bogota each day, how many buses would we have to sample to find the average journey length with 95/10 confidence/precision?

159. The equation used to give us the required number of clusters, c , to sample is:

$$c \geq \frac{1.645^2 MV}{(M-1) \times 0.1^2 + 1.645^2 V} \quad (69)$$

Where:

$$V = \left(\frac{SD}{Cluster\ mean} \right)^2$$

C Number of clusters (buses) to be sampled

M Total number of clusters (buses)

1.96 Represents the 95% confidence required

0.1 Required precision (the equation takes into account that this is relative)

160. To perform the calculations we need information about journey length at the bus level, i.e. total journey length aggregated across all passengers on a bus. If such information does not already exist, we might collect it in a pilot study. The example here assumes that data are available from four buses.

Bus	Total journey length (on average) ⁹
A	195
B	96
C	63
D	159

161. Calculating the mean and standard deviation for this total journey length for a bus gives us:

$$Cluster\ mean\ (i.e.\ \bar{y}) = \frac{1}{n} \sum_{i=1}^n y_i = \frac{195 + 96 + 63 + 159}{4} = 128.25 \quad (70)$$

⁹ These totals may be derived from collecting data on all individual passengers on a bus, or otherwise by taking a sample of them and scaling up from the sample to all the passengers on the bus.

DRAFT

$$SD = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2} = \sqrt{\frac{(196 - 128.25)^2 + \dots + (159 - 128.25)^2}{3}} = 59.7180 \quad (71)$$

162. These statistics (i.e. mean and SD of data) are easily produced using standard statistical software. Substituting these values into the equation gives the required number of clusters, i.e. buses as:

$$V = \left(\frac{59.7180}{128.25} \right)^2 = 0.22 \quad (72)$$

$$c \geq \frac{1.96^2 \times 12000 \times 0.22}{(12000 - 1) \times 0.1^2 + 1.96^2 \times 0.22} = 82.7 \quad (73)$$

163. The total number of buses that should be sampled is 83. Asking the journey lengths from everyone on each of the 83 buses sampled will satisfy the 95/10 confidence/precision criterion.

Example 18 – Multi-stage Sampling

164. Continuing the previous example, suppose that we want to sample a number of local buses, but we only want to sample a number of individuals on each bus (unlike cluster sampling where everyone on each selected bus is sampled). This is an example of multi-stage sampling, as we are sampling a number of groups (buses), and we are then going on to sample a number of units (passengers) from each selected group (bus).

165. We start by assuming that we want to sample five passengers on each selected bus. In general terms we will call this number u (for units).

166. In order to be able to perform a sample size calculation we need information on:

- (a) The variation between individual passengers on the bus;
- (b) The variation between buses;
- (c) The average journey length for a passenger;
- (d) The average journey length at the bus level (when aggregated across all passengers).

167. A previous study had provided data for passengers on three different buses, and the results are summarized below.

168. Note that not all buses will have exactly the same number of passengers.

DRAFT

Journey length (km)				
Bus	Number of passengers	Mean journey length ¹⁰	Total journey length aggregated over all passengers on the bus	Standard deviation ¹¹ (between passengers on the same bus)
A	26	6.9	179	3.30
B	21	7.5	157	6.21
C	30	6.7	200	3.78
Total number of passengers	77			
Overall mean journey length per passenger		7.0		
Mean total journey length (per bus)			179	
SD_B = Standard deviation between buses (SD of the total journey length column)			4889	
SD_W = Average between passenger (within bus) standard deviation				4.44

169. In the table above, the overall mean journey length is the average length per passenger, i.e.

$$\text{Overall mean} = \frac{179 + 157 + 200}{77} = 7.0 \quad (74)$$

170. The cluster mean journey length is the mean length per bus, i.e.

$$\text{Cluster mean} = \frac{179 + 157 + 200}{3} = 179 \quad (75)$$

171. SD_W^2 is the average variance between passengers within buses. Its square root (SD_W) is the average within bus standard deviation. The equation for SD_W^2 is:

$$SD_W^2 = \frac{26 \times 3.30^2 + 21 \times 6.21 + 30 \times 3.78^2}{77} = 19.75 \quad \text{and so} \quad SD_W = 4.44 \quad (76)$$

172. SD_B^2 is the variance between the mean journey lengths per bus. Its square root is the standard deviation between buses. It can be calculated using the general equation for a variance:

$$SD_B^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1} \quad \text{where the } y_i \text{ are the journey lengths for the different buses.}$$

$$SD_B^2 = 467 \quad \text{and so} \quad SD_B = 22 \quad (77)$$

¹⁰ This can be a mean from all passengers on the bus or a mean from a sample of passengers.

¹¹ And this can be a standard deviation based on all passengers or from a sample of passengers.

DRAFT

173. As well as the information from the table above we also require the number of buses, and the average number of passengers on each bus for the whole population. For this example we are using 12,000 buses with an average of 30 passengers on each bus.

$$c \geq \frac{\left(\frac{SD_B}{Clustermean}\right)^2 \times \left(\frac{M}{M-1}\right) + \left(\frac{1}{u}\right) \times \left(\frac{SD_W}{Overallmean}\right)^2 \left(\frac{\bar{N}-u}{\bar{N}-1}\right)}{\left(\frac{0.1}{1.645}\right)^2 + \frac{1}{M-1} \left(\frac{SD_B}{Clustermean}\right)^2} \quad (78)$$

Where:

M	Total number of groups (12,000 buses)
\bar{N}	Average number of units per group (30 passengers per bus)
u	Number of units that have been pre-specified to be sampled per group (pre-specified number of passengers to be sampled on each bus = 5)
1.96	Represents the 95% confidence required
0.1	Required precision

$$c \geq \frac{\left(\frac{22}{179}\right)^2 \times \left(\frac{12000}{12000-1}\right) + \left(\frac{1}{5}\right) \times \left(\frac{4.44}{7.0}\right)^2 \left(\frac{30-5}{30-1}\right)}{\left(\frac{0.1}{1.96}\right)^2 + \frac{1}{12000-1} \left(\frac{22}{179}\right)^2} = 32.6 \quad (79)$$

174. Therefore if we were to sample five passengers from each bus we should sample 33 buses for the required confidence/precision.

175. Producing a table such as the one below, with different values of u , can help decide the practicalities of allocating limited resources, while still satisfying the 95/10 confidence/precision criterion.

Number of passengers sampled on each bus	Required number of buses
u	c
5	33
10	17
15	12
20	9
25	7

176. In this example, by doubling the number of passengers on each bus to be sampled from 5 to 10, we substantially reduce the number of buses that need to be sampled from 33 to 17.

177. Note that in the above example the numbers of passengers on a bus were different for the different buses. In practice this is likely to be the case, although the actual numbers may not always be known. This is not critical to the sample size calculation. What is important is that sensible estimates of the mean and standard deviation at both the cluster level (bus level) and unit level (passenger level) are used in the calculation.

DRAFT

Mean value parameter of interest (Transport Project)

178. This section covers an example sample size calculation based on systematic sampling where the objective of the project relates to a mean value of interest.

179. As for all absolute parameter of interest examples we need to know:

- (a) The expected mean (the desired reliability is expressed in relative terms to the mean);
- (b) The standard deviation;
- (c) The level of precision, and confidence in that precision (95/10 for all large-scale examples).

180. The parameter of interest for the example below is the average journey time of buses on a specific route. We know that over a month 960 journeys are made on the route of interest, and the average journey time is 18 minutes with a standard deviation of 6 minutes.

Example 19 – Systematic Sampling

181. Using systematic sampling we want to sample every n^{th} journey on that route. The sample size equation for a required 95/10 confidence/precision is:

$$n \geq \frac{1.96^2 V}{0.1^2} \quad (80)$$

Where:

$$V = \left(\frac{SD}{mean} \right)^2$$

182. Substituting in mean and standard deviation from above gives:

$$V = \left(\frac{6}{18} \right)^2 = 0.11 \quad (81)$$

$$n \geq \frac{1.96^2 \times 0.11}{0.1^2} = 42.7 \quad (82)$$

183. In total we should sample 43 journey times. We want to take these samples evenly spread over the 960 journeys made each month, therefore we should sample one journey for every N/n – that is one journey every 22 (= 960 / 43).

184. We can make sure that we sample at random by selecting a random starting point between 1 and 22, say journey 18, and then sample every 22nd journey from this point on: 18, 40, 62, 84, 106, etc. up to 960. This would give us a sample evenly spread over the month that is large enough to estimate the average journey time with 95/10 confidence/precision.

185. It may be more practical to sample every 20th journey rather than every 22nd. This would result in more samples being taken than the 43 calculated above – the only effect this would have would be to increase the precision and so would be perfectly acceptable.

DRAFT

Sampling a proportional parameter of interest when N is small or p is very low or very high

186. The following is an example of how to deal with very small or very large proportions in the context of simple random sampling.

187. If N (the population total) is very small, or p (the proportion of interest) is very low or very high then the methods of calculating sample sizes described above are not appropriate. This is because the normal approximation of the number of events (Np) will not be suitable, and this is an underlying assumption as the sample size equations are based on normal confidence intervals. It is recommended that the equations illustrated previously in this document are only used when Np and $N(1-p)$ are both greater than 10. When this is not the case the equation below should be used:

$$n = -z^2(1-2\theta) + z^2 \sqrt{\frac{1}{(2 \times \text{precision} \times p)^2} - 4\theta + 4\theta^2} \quad (83)$$

Where:

$$\theta = \frac{p(1-p)}{(2 \times \text{precision} \times p)^2}$$

n Sample size required

z Z value relating to the level of confidence (95% = 1.96, 90% = 1.645)

p Estimate of the proportion of interest

precision Relative precision

188. For example, if the proportion of interest is 0.03, with 95%/10% confidence and precision then

$$\theta = \frac{0.03(1-0.03)}{(2 \times 0.1 \times 0.03)^2} = 808 \quad (84)$$

$$n = -1.96^2(1-2 \times 808) + 1.96^2 \sqrt{\frac{1}{(2 \times 0.1 \times 0.03)^2} - 4 \times 808 + 4 \times 808^2} = 12447 \quad (85)$$

189. Once the sampling has been completed and the sample proportion calculated (p^{\wedge}) then a 95% confidence interval should be calculated as follows:

$$\frac{A-B}{C} \text{ to } \frac{A+B}{C} \quad (86)$$

Where:

$$A = 2np^{\wedge} + 1.96^2$$

$$B = 1.96 \sqrt{1.96^2 + 4np^{\wedge}(1-p^{\wedge})}$$

$$C = 2(n + 1.96^2)$$

DRAFT

(Replace 1.96 with 1.645 for a 90% confidence interval).

II. Reliability calculations

A. Introductory notes

190. The two examples presented here illustrate how to estimate a numeric parameter and a proportion and how to check their reliability. The sampling method used in both cases is simple random sampling. Both examples are assumed to be small-scale project activities where the required reliability criteria is 90/10, i.e. 90% confidence and 10% precision.

191. If calculations are being performed manually, it is important to retain as many decimal places as relevant, until the final calculated figure is reached. Rounding can then be carried out. To emphasize this, the calculations presented here use figures to several decimal places.

Example 1: CFL Project – Numeric parameter

192. The parameter of interest in this example is the mean average daily usage of a CFL (in hours) for a whole population of CFLs that were distributed in a particular region of a country.

193. The population is the 420,000 households to which CFLs were distributed, one per household. A simple random sample of 140 households was taken, and the average daily usage (in hours) of each CFL was recorded. These are presented in the table below.

Average CFL usage (in hours)

cfl	usage	cfl	usage	cfl	usage	cfl	usage	cfl	usage	cfl	usage	cfl	usage
1	3.78	21	3.63	41	2.81	61	4.17	81	3.62	101	2.24	121	0.58
2	3.12	22	3.17	42	4.57	62	4.68	82	2.46	102	4.79	122	6.09
3	4.42	23	3.26	43	3.56	63	2.99	83	6.14	103	4.59	123	0.39
4	4.09	24	6.97	44	4.41	64	3.34	84	0.67	104	3.27	124	3.69
5	1.15	25	0.48	45	3.26	65	5.37	85	4.73	105	1.86	125	2.04
6	2.87	26	2.50	46	0.30	66	2.17	86	1.03	106	0.00	126	4.51
7	4.79	27	2.92	47	5.48	67	2.36	87	2.34	107	6.70	127	4.39
8	4.20	28	6.82	48	1.75	68	3.12	88	4.66	108	3.36	128	3.58
9	1.13	29	0.92	49	3.38	69	4.69	89	2.40	109	5.39	129	4.23
10	3.68	30	2.35	50	1.24	70	5.40	90	5.28	110	2.04	130	5.28
11	2.91	31	0.19	51	3.62	71	4.22	91	5.90	111	3.58	131	3.71
12	2.47	32	4.19	52	7.41	72	1.27	92	0.60	112	6.27	132	2.41
13	3.46	33	3.15	53	1.74	73	2.93	93	5.85	113	0.41	133	1.58
14	2.19	34	3.19	54	3.60	74	2.17	94	1.22	114	4.55	134	3.96
15	2.25	35	7.15	55	2.18	75	4.24	95	7.76	115	2.61	135	5.86
16	2.37	36	1.70	56	4.12	76	6.07	96	4.50	116	6.37	136	5.46
17	2.38	37	2.98	57	4.88	77	5.26	97	5.68	117	4.30	137	2.90
18	3.23	38	5.00	58	2.92	78	2.46	98	2.81	118	3.08	138	3.17
19	1.78	39	0.99	59	0.82	79	1.33	99	4.03	119	3.17	139	4.17
20	3.57	40	6.54	60	3.16	80	2.55	100	0.24	120	6.24	140	6.93

194. The parameter of interest – the mean average daily usage of a CFL (in hours) for this whole population of CFLs – is estimated from the sample mean. This is often written as \bar{y} (and is

equal to $\frac{1}{n}(y_1 + y_2 + \dots + y_n)$, or the shorthand form $\frac{1}{n} \sum_{i=1}^n y_i$). n is the sample size, i.e. 140.

**DRAFT**

195. The mean average usage for the sample of 140 CFLs is 3.4686 hours. As a simple summary this is rounded to 1 or 2 decimal places, i.e. the mean average usage of the CFLs is estimated to be 3.47 hours.

Confidence, precision and reliability

196. Instead of presenting just a single estimate, it is better to summarize the results of sampling using a confidence interval. In this example the 90% confidence interval is 3.22 to 3.71 hours. We are 90% sure that the true population mean value for average usage of a CFL is between 3.22 hours and 3.71 hours.

197. The 90% confidence interval for the population mean is given by the equation: sample mean \pm t-value \times standard error of the mean.

198. The estimate of 3.47 hours is regarded as reliable if the precision of the study – as defined by the t-value \times standard error of the mean – is within the pre-specified reliability precision. For small-scale mechanisms this is 10% of the mean.

199. Detailed calculations are presented below. In this example the precision is 7.1% of the mean and so the sample estimate of 3.47 hours is within the required specification.

Checking reliability**(i) Standard error of the mean**

200. The equation for the standard error of the mean when data have been collected using

simple random sampling is $\sqrt{(1-f) \frac{s^2}{n}}$.

f is the sampling fraction – the proportion of the population that is sampled.

Here it is $\frac{140}{420000} = 0.0003$.

s^2 is the sample variance (s is the sample standard deviation).

For this sample of 140 CFLs, $s^2 = 3.0826$ and $s = 1.7557$.

n is the sample size, i.e. 140.

201. Putting all these pieces of information together gives:

$$\sqrt{(1-f) \frac{s^2}{n}} = \sqrt{\left(1 - \frac{140}{420000}\right) \times \frac{3.0826}{140}} = \sqrt{0.99967 \times \frac{3.0826}{140}} = \sqrt{0.0220} = 0.1484 \quad (87)$$

and so the standard error of the mean is 0.1484.

DRAFT

(ii) t-value

202. This value depends on (i) the level of confidence and (ii) the size of the sample. The exact figure can be acquired from statistical tables for the t-distribution, or using standard statistical software. The value can also be derived in Microsoft Excel using the TINV¹² function.

For a sample size of 140 the t-value is 1.6559.

(iii) Precision

203. The precision associated with an estimate is: t-value × standard error of the mean.

The precision of the mean average CFL usage (in hours), assuming 90% confidence, in this example is therefore: ± (1.6559 × 0.1484) i.e. ± 0.2457.

The ratio of this relative to the mean CFL usage is $\frac{0.2457}{3.4686} = 0.0708$ and so the relative precision is 7.1%. The data are therefore within the required specification.

Another way of checking reliability

204. The limits of the confidence interval are sample mean ± t-value × standard error of the mean, which can be written more generally as sample mean ± precision, where the lower limit is mean minus precision and the upper limit is mean + precision.

205. Reliability can therefore be checked using the following calculation:

$$\frac{\frac{1}{2} \text{width of confidence interval}}{\text{mean}} \times 100\% \quad (88)$$

206. For example, here the mean CFL usage is 3.4686, and the 90% confidence interval is 3.2230 to 3.7143 hours. Reliability is therefore:

$$\frac{\frac{1}{2}(3.7143-3.2230)}{3.4686} \times 100\% = \frac{\frac{1}{2} \times 0.4913}{3.4686} \times 100\% = 7.1\% \quad (89)$$

207. The above approach is likely to be most useful when the data have been analysed using statistical software which produced the relevant confidence interval as well as the sample mean.

Example 2 : Cook Stove Project – Proportional parameter

208. The parameter of interest in this example is the proportion (or percentage) of cook stoves in a particular region of a country that were still operational at the end of the third year after the stoves were distributed.

209. The population of interest is the 640,000 households, and there was one cook stove per household. A simple random sample of 274 of these households was taken, and for each of them it was recorded whether or not the cook stove was still operational.

¹² TINV(0.10,(sample size minus 1)) will give the t-value associated with 90% confidence. For example here TINV(0.10,139) gives the t-value for a sample size of 140 and 90% confidence.

DRAFT

210. The parameter of interest – the proportion (or percentage) of cook stoves that were still operational in the whole population – is estimated from the sample proportion.

211. This is often written as p and is calculated as $p = \frac{r}{n}$ where r is the number of “successes”, in this case the number of cook stoves that are still in operation, and n is the total number of cook stoves that are observed in the sample.

212. In this example there were 159 cook stoves out of the 274 that were still in operation. The sample proportion is therefore $p = \frac{159}{274} = 0.5803$. Rounding this to two decimal places gives us a proportion of 0.58. In other words, 58% of the cook stoves were still operational after the third year.

Confidence, precision and reliability

213. Instead of presenting just a single estimate, it is better to summarize the results of sampling using a confidence interval. In this example the 90% confidence interval for the proportion is 0.5313 to 0.6293. We are therefore 90% sure that the percentage of cook stoves in the population that are still operational is between 53% and 63%.

214. The 90% confidence interval for the population proportion is given by the equation: sample proportion $\pm 1.6449 \times$ standard error of the proportion.¹³

215. The estimate of 58% is regarded as reliable if the precision of the study – as defined by $1.6449 \times$ standard error of the proportion – is within the pre-specified reliability precision. For small-scale mechanisms this is 10% of the proportion. In this case ± 0.058 in absolute terms or $\pm 5.8\%$.

216. Detailed calculations are presented below. In this example the precision is 8.5% of the sample proportion and so the sample estimate of 58% operational cook stoves is within the required specification.

Checking reliability**(ii) Standard error of the proportion**

217. The equation for the standard error of the proportion when data have been collected using simple random sampling is $\sqrt{(1-f) \frac{pq}{n}}$

f is the sampling fraction – the proportion of the population that is sampled.

Here it is $\frac{274}{640000} = 0.00043$

p is the sample proportion, i.e. 0.5803. $q = (1-p)$. It represents the proportion of cook stoves that are not operational after three years, and is 0.4197. n is the sample size, i.e. 274.

¹³ A confidence interval for a proportion is: sample proportion \pm z-value \times standard error of the proportion. The z-value depends on the level of confidence. For 90% confidence it is 1.6449.

DRAFT

218. Putting all these pieces of information together gives:

$$\sqrt{(1-f)\frac{pq}{n}} = \sqrt{(1-0.00043)\frac{0.5803 \times 0.4197}{274}} = \sqrt{0.00089} = 0.0298 \quad (90)$$

219. Note that this standard error could also be calculated using the actual numbers of population size, sample size, number of operational cook stoves etc., i.e.:

$$\sqrt{(1-f)\frac{pq}{n}} = \sqrt{\left(\frac{640000-274}{640000}\right)\frac{\left(\frac{159}{274}\right)\left(\frac{115}{274}\right)}{274}} = \sqrt{0.00089} = 0.0298 \quad (91)$$

220. The standard error of the proportion is 0.0298. In terms of the standard error of the percentage it is 2.98%.

(ii) Precision

221. The precision associated with a proportion is: z-value × standard error of the proportion. The precision of the proportion of operational cook stoves in this example, assuming 90% confidence, is: ± (1.6449 × 0.0298) i.e. ± 0.0490.

222. The ratio of this relative to the proportion of cook stoves that are still operational is $\frac{0.0490}{0.5803} = 0.0845$ and so the relative precision is 8.5%. The data are within the required specification.

Another way of checking reliability

223. The limits of the confidence interval are sample proportion ± z-value × standard error of the proportion, which can be written more generally as sample proportion ± precision, where the lower limit is the proportion minus precision and the upper limit is the proportion plus precision.

224. Reliability can therefore be checked using the following calculation:

$$\frac{\frac{1}{2}\text{width of confidence interval}}{\text{proportion}} \times 100\% \quad (92)$$

For example, here the proportion of cook stoves that are still operational is 0.5803, with a 90% confidence interval of 0.5313 to 0.6293.

Reliability is therefore:

$$\frac{\frac{1}{2}(0.6293 - 0.5313)}{0.5803} \times 100\% = \frac{\frac{1}{2} \times 0.0980}{0.5803} \times 100 = 8.5\% \quad (93)$$

225. The above approach is likely to be most useful when the data have been analysed using statistical software which produced the relevant confidence interval and sample proportion.

Comments

**DRAFT**

226. The equation above assumes that the distribution of the proportion is approximately Normal. That is usually an acceptable assumption provided the proportion of interest is not too small and not too large, and the sample size is not too small;

227. If the sampling fraction f is small then the multiplier $(1 - f)$ in the above calculation will be very close to 1. In some instances, therefore, the equation that is used for the standard error of the proportion is the conservative equation $\sqrt{\frac{pq}{n}}$ ¹⁴;

228. If statistical software is used to undertake the calculation, the software may use the exact equation for calculating the confidence interval (which assumes a Binomial distribution as opposed to the Normal approximation). In this case the reliability would be checked using the equation which is based on the width of the confidence interval.

¹⁴ It is conservative because $\sqrt{\frac{pq}{n}}$ will be greater than $\sqrt{(1 - f)\frac{pq}{n}}$.



DRAFT

Appendix

Sampling Scheme	Advantages	Disadvantages
Simple Random Sampling: Taking a random sample from the whole population.	Easiest method to understand and therefore use. Suitable if there is little heterogeneity amongst the units being sampled.	Requires knowledge of entire population before a sample can be selected. If the population covers a large geographical area, then it can often lead to sampling units that are spread out over the area. Such a situation can often be costly. Only suitable if the population being studied is relatively homogeneous with respect to the parameter being studied.
Systematic Sampling: Taking a sample every n units.	Easy to apply. Commonly used as it ensures there is always sufficient distance between samples.	Leads to units being spread out over a large geographic area. Such a geographic distribution can often be costly.
Stratified Random Sampling: Randomly sampling a different number of units from each strata according to the weight of each strata in the population.	Improves the precision of the estimate (compared to simple random sampling) if there are differences between the strata.	Complicated to calculate. What the stratification factors should be is not always obvious.
Cluster Sampling: Sampling every unit in a sample of n clusters from the population.	The most economical form of sampling as units are all grouped according to one criterion (often geographical). Sometimes the only approach, since a list of all households may not be available, only a list of villages. Once the villages have been selected, the households can be sampled. It saves time at a management level.	Results are not normally so 'good' (i.e. standard errors of estimates tend to be high due to homogeneity of characteristics in the subgroup sampled). [But a larger sample can help to compensate for this]

**DRAFT**

Sampling Scheme	Advantages	Disadvantages
Multi-stage Sampling: Randomly sampling a number of units within a number of randomly selected clusters.	Enables sampling approach at two levels. Can compare different scenarios – number of clusters and number of units within the clusters – in order to find most cost-efficient and reliable scenario.	Analysis and the sample size calculation are more difficult.

History of the document

Version	Date	Nature of revision
01.0	EB 66, Annex # 02 March 2012	Initial adoption.
Decision Class: Regulatory Document Type: Guideline Business Function: Methodology		