## FACTSHEET III

If you have decided to use a survey approach for obtaining customer feedback, you need to determine what sample size to use.  This Factsheet first discusses sample sizes, sampling error, and confidence intervals—all of which factor into decisions about the sample size.  It then presents a table for you to use in determining what sample size to use—and tells you step by step how to make use of that table.  This Factsheet then describes how to go about randomly selecting that number of customers from the total list of customers you have served during the time period to be covered by the survey.

### WHAT KINDS OF SAMPLE SIZES ARE WE TALKING ABOUT?

Before we give specific guidelines on how to choose the sample size, it will be useful to set some general expectations.  National public opinion polls like the Gallup Poll and the Roper Poll typically use sample sizes in the range of 1,350 to 1,800.  These polls use fairly large sample sizes to obtain a result that represents the entire adult U.S. population with a sampling error on the order of plus or minus 2.5 percent to 3 percent.  Such small levels of sampling error are needed because the polls often address matters of national importance.  The decisions made, based in part on the results of these national polls, may be far-reaching, long-lasting, and affect millions of people.

The surveys you will be conducting to obtain customer feedback will be of a very different nature.  The target group whose opinions you need will be much smaller: It will probably be the people who have come to you and your colleagues in one specific program area within EPA, within a limited time (e.g., during one year) to request certain products or services.  We are therefore talking about a target group of maybe as many as 500 to 1,000 people (few EPA programs directly serve more customers than that) and in some cases 50 people or fewer.  Furthermore, although the decisions that will be affected by customer feedback are important, they will probably not be far-reaching and long-lasting.  The scope of decisions to be made in most cases will be, for example:

- Should we change a process to reflect customer comments?
- Should we revise some of our written products?
- Should we provide a half day of customer feedback training to each staff member?

Even in the worst case—we make the wrong decision about whether our products need to be revised and whether the staff members need further training—we will (if we continue to obtain feedback from our customers at least once each year) discover our error soon enough and be able to correct it, without incurring excessive or irreparable damage in the meantime.

Based on these considerations, it is reasonable to have higher sampling errors than those associated with national surveys like the Gallup Poll. We can feel comfortable with sampling errors of 5 percent or even 10 percent.

Additionally, for getting feedback from EPA customers, we have relatively small target groups who were served by a specific program during the time period of immediate interest. For this reason, it is reasonable for you to use a much smaller sample size than is used in the Gallup and Roper Polls, which seek to accurately capture the opinions of millions of people.

### *Sampling error*

"Sampling error" is normally presented as a percentage with a plus or minus sign in front of it. For example, the sampling error in one particular situation may be ± 3.5 percent. That means that the *true* value of a given measure for the entire population—that is, the whole target group you are getting feedback from—is the value obtained from your sample of customers, plus or minus 3.5 percent. If for example, 62.4 percent of your sampled customers are satisfied, the actual percentage of satisfied customers lie within the range between 58.9 percent (62.4 percent - 3.5 percent) and 65.9 percent (62.4 percent + 3.5 percent).

But that is not quite true. In fact, there is no range of reasonable size that we can identify for which we can be certain that the true value for the full list of customers lies in that range.

Why is that so?

Because there's always the possibility of very unlikely circumstances occurring—with the result that the characteristics of the customers in the sample are very different from the characteristics of the customers not in the sample. In such circumstances, the true value for all customers will be very different from the value obtained from the customers in the sample surveyed. The only way to get around this statistical fact is to specify "how certain we want to be" that the true value does, in fact, fall with a specific range around the value we obtain from the sample. This degree of certainly we are looking for is known as the "confidence level."

### *Confidence level*

The confidence level indicates how confident we want to be that the true value lies within a specific range.

There is no one confidence level that is the right one to use. There are many different possible confidence levels, and only you can decide which confidence level is appropriate for your survey.

Much of the work in the area of public opinion surveys uses the 95 percent confidence level. That means that if you determine the sampling error using the 95 percent confidence level, you can be

95 percent certain that the true value for all your customers will lie within a specific percentage band (one equal to the size of the sampling error) around the result you obtain from the sample of customers you contact.

Another confidence level commonly used is the 90 percent confidence level. With a 90 percent confidence interval, you can be confident that 9 times out of 10, the true value falls within the value obtained from your sample of customers, plus or minus the sampling error. Some analysts use 80 percent confidence intervals.

To decide what confidence level to use, you might want to think of a scale running from 80 to 95, where 95 represents a high level of confidence and 80 represents a lower level of confidence. Decide which confidence level to use based on the way in which your results will be used, how products and services may be affected by the results, and the frequency with which you will collect additional information to confirm or revise your findings.

### *Determining the sample size*

Now that we have established appropriate expectations with regard to sampling error and sample size, we will provide you with some guidance on selecting your sample size. Please recognize that there are several factors to consider in determining the sample size. The information provided here is intended to help get you started. Please refer as well to the additional information provided in **Factsheets III, IV,** and **V**. If you wish, you may also consult a statistician within your Office at EPA. A list of EPA statisticians showing the EPA Office in which each of these statisticians is located can be obtained from the Office of the Chief Statistician of EPA within EPA's Center for Environmental Information and Statistics by calling 202-260-5244.

| Number in target group | Sampling error | Confidence level | Sample size |
|---|---|---|---|
| 1000 | ±5 | 80 | 141 |
| 1000 | ±5 | 90 | 214 |
| 1000 | ±5 | 95 | 278 |
| 500 | ±5 | 80 | 124 |
| 500 | ±5 | 90 | 176 |
| 500 | ±5 | 95 | 218 |
| 200 | ±5 | 80 | 90 |
| 200 | ±5 | 90 | 116 |
| 200 | ±5 | 95 | 132 |

| Number in target group | Sampling error | Confidence level | Sample size |
|:---:|:---:|:---:|:---:|
| 100 | ±5 | 80 | 62 |
| 100 | ±5 | 90 | 74 |
| 100 | ±5 | 95 | 80 |
| 50 | ±5 | 80 | 39 |
| 50 | ±5 | 90 | 43 |
| 50 | ±5 | 95 | 45 |
| 1000 | ±10 | 80 | 39 |
| 1000 | ±10 | 90 | 64 |
| 1000 | ±10 | 95 | 88 |
| 500 | ±10 | 80 | 38 |
| 500 | ±10 | 90 | 60 |
| 500 | ±10 | 95 | 81 |
| 200 | ±10 | 80 | 34 |
| 200 | ±10 | 90 | 51 |
| 200 | ±10 | 95 | 66 |
| 100 | ±10 | 80 | 29 |
| 100 | ±10 | 90 | 41 |
| 100 | ±10 | 95 | 50 |
| 50 | ±10 | 80 | 23 |
| 50 | ±10 | 90 | 29 |
| 50 | ±10 | 95 | 34 |

The above table is that appropriate for *simple random sampling* (SRS), which is a sampling procedure based on sampling without replacement. Simple random sampling is the most commonly used sampling procedure. The table is based on the *approximate formula* given in **Factsheet IV**. This approximate formula includes an adjustment comparable to the finite population correction factor for each combination of target population and sample size.

The *precise formula* that can be used instead of this *approximate formula* is also given in **Factsheet IV**. For a discussion of the *finite population correction factor*, see **Factsheet IV**. For

a discussion of the meaning and significance of *sampling without replacement* (as contrasted with *sampling with replacement*), see the discussion of this matter in the last section of **Factsheet V**.

The procedure described below in this Factsheet for randomly selecting a sample from the full list of customers served in a specific period of time is *simple random sampling* and is therefore consistent with the above table.

### *Here's how to use the above table*

The instructions that follow assume that the unit of analysis for the survey will be the "person served." (See **Factsheet VII** for a discussion of "Unit of Analysis.")

1) Identify the number of persons you have served in the time period of interest. Find that number in the column labeled "Number in Target Group."

2) Select the confidence level that you consider to be the most appropriate given the magnitude of the decisions that will be made based (in part) on the results obtained from the survey:

   • If the decisions to be made using the survey results will be far-reaching, long-lasting and/or costly, use the 95 percent confident level

   • If the decisions to be made using the survey results will be less far-reaching, less long-lasting or less costly, use the 90 percent confidence level

   • If the decisions to be made using the survey results will have more limited consequences, mostly in the short-term (e.g., in the next 6 to 12 months) and the cost implications of the decisions will be moderate, you may use the 80 percent confidence level.

3) Select the level of sampling error you consider to be acceptable given the magnitude of the decisions that will be made using the results obtained from the sample.

   • For most EPA customer satisfaction surveys, a sampling error of ±10 percent should be acceptable.

   • In cases where the decision to be made based (in part) on the survey results is of such a nature that a smaller level of sampling error is needed, a sampling error of ±5 percent can be used instead.

4) Read off the corresponding sample size.

- If the total number of customers served falls between two of the values shown above in the column "Number in Target Group," you can use interpolation to obtain an initial estimate of the appropriate sample size.

- You can then use the approximate formula for determining sample size presented in **Factsheet IV** to obtain a much better estimate of the sample size needed.

- You can stop here and make use of *the approximate value for the sample size* obtained in step 4) b) immediately above.  Alternatively, you can, if you wish, now make use of *the trial and error approach* presented in **Factsheet IV** or, even better, *the combined approach*, also presented in **Factsheet IV**, to calculate *the precise value for the sample size* needed.

### Here's how to randomly select a sample of customers once you have determined what sample size to use

Once you have determined the appropriate sample size to use, the next step is to randomly select that number of customers from the total number served in the time period of interest.  Here is a procedure you can use to make that random selection:

1) Make a complete list of all the persons served in the period of interest for which you already have (or can obtain, with a reasonable expenditure of effort) the needed contact information (i.e., name, plus address or phone number).  Put the customers in alphabetical order to ensure that there are no duplicate names.  Eliminate any duplicate names.

2) Once all duplicate names have been eliminated (so that each name appears only once), starting at the top of the list, number each name.  The result is the *master list* of customers served. The number next to each name is that person's customer number.

3) Here is a computer based approach for selecting a sample of customers from the *master list*:

   - You will use spreadsheet software (like Lotus 1-2-3 or Microsoft Excel) to carry out the remaining steps of this procedure.  Before you begin to make use of any particular spreadsheet software, first make sure that it has a "randomize" function.  Not all spreadsheets do.

   - Enter the customer numbers in numerical order into the spreadsheet, one number per row. Place each of these numbers in the second column of the spreadsheet, leaving the first column in each row blank. The result will be a spreadsheet with the number of rows equal to the number of customers and with the rows having the numbers 1, 2, 3, and so on (up to the total number of customers served), with these numbers in the second column of each row.

- Use the randomize function on the second column of the spreadsheet. The numbers in the second column are now in random order.

- Enter numbers into the first column of each row. Enter the number 1 into this column in the first row, enter 2 into this column in the second row, and so on. These new numbers are the row labels.

- Mark off the number of rows corresponding to the sample size chosen above. For example, if the sample size is 65, mark off the first 65 rows.

- The numbers appearing in the second column of the rows marked off in step e) above are the customer numbers corresponding to the customers to be included in the sample. For each of these customer numbers, read off the name of the customer appearing next to this number on the master list prepared in step 2) above and place it in a new list. This is new list is the list of customers selected for inclusion in the sample—the people you will contact during the survey and ask to respond to the survey questions.

- If due to a lower than expected response rate, the number of customers from whom responses are received is less than the desired sample size, and all reasonable followup efforts have already been made to increase the response rate, go back to the spreadsheet and mark off the additional number of rows needed to reach the desired sample size. The numbers appearing in the second column of these additional rows are the customer numbers for the additional customers to be added to the sample.

For an equivalent procedure that does not make use of a computer or a computer spreadsheet, see the last section of **Factsheet V**.