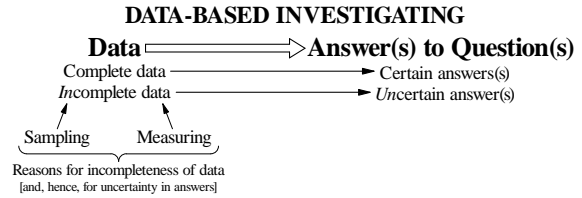
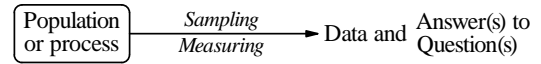


**Figure 6.1. FROM PROBABILITY TO STATISTICS: Modelling Sample Error and Measurement Error**

As summarized in the two schemas at the right, statistics is concerned with *data-based investigating* of some population or process to *answer* one or more (statistical) *questions* of interest.

- \* If the investigating yields *complete* information, we can obtain a *certain* answer; that is, an answer we can *know* is correct.
- \* If the investigating yields *incomplete* information, we *cannot* know an answer is correct (an *uncertain* answer) – in fact, it is *unlikely* a *numerical* answer (an ‘estimate’) is (exactly) correct;
  - sampling and measuring yield data (and, hence, information) that are *inherently* incomplete.



In this Figure 6.1, we develop a probability a model for the investigative processes of sampling (which involves selecting and estimating) and measuring; this model will allow us, in Figure 6.3, to quantify the *likely* size of sample error and measurement error – that is, to quantify uncertainty from these two sources – for numerical answers to some types of questions.

**1. Ideas from Probability**

To go from probability to statistics, ideas we need from probability modelling in the previous Part 5 are as follows.

- \* Using the normal distribution to model the shape of appropriate data distributions, as summarized by the probability statement (6.1.1) about the random variable  $Y$  at the right.  $Y \sim N(\mu, \sigma)$  ----(6.1.1)
- \* A random variable which is a sum or an average of (probabilistically independent) normal random variables also has a normal distribution and its mean and standard deviation can be expressed in terms of the mean(s) and standard deviation(s) of the random variables that make up the sum or average. For example, for  $n$  probabilistically independent  $N(\mu, \sigma)$  random variables  $Y_j, j=1, 2, \dots, n$ , as given in equations (6.1.2) and (6.1.3) at the right below:
  - their sum ( $T$ ) [or *total*] has a normal distribution with a *larger* mean and standard deviation than the individual random variables by respective factors of  $n$  and  $\sqrt{n}$ .  $T \sim N(n\mu, \sigma\sqrt{n})$  ----(6.1.2)
  - their *average* ( $\bar{Y}$ ) has a normal distribution with the *same* mean  $\mu$  as the individual random variables but a standard deviation that is *smaller* by a factor of  $1/\sqrt{n}$ .  $\bar{Y} \sim N(\mu, \sigma\sqrt{1/n})$  ----(6.1.3)

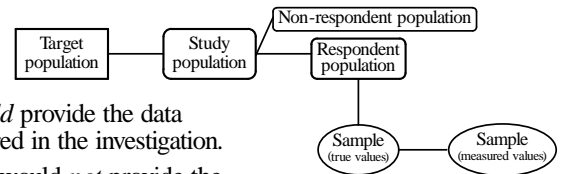
If the requirement for *normality* of the distribution of the individual random variables in a total or an average is relaxed, the Central Limit Theorem (CLT) approximation allows us to write equations (6.1.4) and (6.1.5); the *accuracy* of this approximate normality of a total or an average improves as:

- o the value of  $n$  gets *larger*;
  - o the distribution(s) of the  $Y$ s get *more symmetrical*.
- $T \doteq N(n\mu, \sigma\sqrt{n})$
- (6.1.4)
- 
- $\bar{Y} \doteq N(\mu, \sigma\sqrt{1/n})$
- (6.1.5)

**2. Ideas from Statistics**

The schema at the right reminds us that data-based investigating in statistics is concerned with five groups of elements or units and their **attributes**: a quantity defined as a function of response (and, perhaps, explanatory) variates over the group (e.g., an average).

- \* **Target population**: the group of elements to which the investigator(s) want Answer(s) to the Question(s) to apply.
- \* **Study population**: a group of elements *available* to an investigation.
- \* **Respondent population**: those elements of the study population that *would* provide the data requested under the incentives for response offered in the investigation.
- \* **Non-respondent population**: those elements of the study population that *would not* provide the data requested under the incentives for response offered in the investigation.
- \* **Sample**: the group of units selected from the respondent population *actually used* in an investigation – the sample is a *subset* of the respondent population (as indicated by the *vertical* line in the schema above).

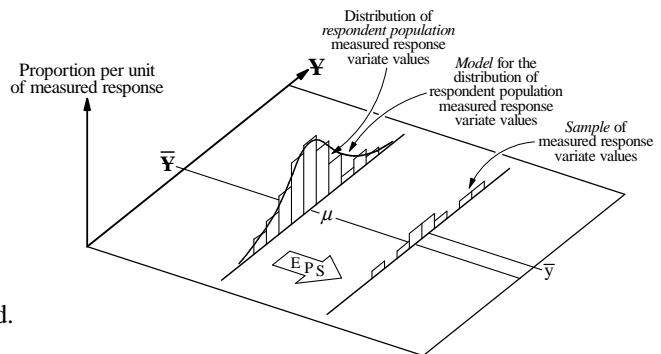


As indicated pictorially in the 3-dimensional diagram at the right, we think of a *respondent* population as having a *distribution* of values for some measured *response* variate ( $Y$ ) of interest; this distribution can be displayed as a *histogram*. Simple population *attributes* of interest are:

- \* Size:  $N$  [the number of *elements*];
- \* Location:  $\bar{Y}$  [the *average*];
- \* Variation:  $S$  [the (*data*) *standard deviation*];

Unless  $N$  is to be estimated, it is usually considered to be:

- o known (or ignored), and
- o large compared to the size ( $n$ ) of any *sample* (to be) selected.



(continued overleaf)

The probability model (6.1.1) for the shape of the histogram is included in the diagram overleaf on page 6.3 at the bottom right, which also reminds us that  $\mu$ , the mean of the model, represents the respondent population average,  $\bar{Y}$ , of the response variate.

The *sample* average is  $\bar{y}$  and equation (6.1.3) is a model for the corresponding random variable  $\bar{Y}$ ; the method of *selecting* the sample [equiprobable selecting (EPS)] is (part of) the reason we can treat the measured response variate values,  $y_j, j=1, 2, \dots, n$ , of the  $n$  units in the sample as values,  $y_j, j=1, 2, \dots, n$ , of random variables,  $Y_j, j=1, 2, \dots, n$ , for the model (6.1.3).

**3. The Set of All Possible Samples from a Population**

Our interpretation of model quantities involves the idea of repeating over and over processes like selecting and measuring; for example, for the random variable  $Y_j$ , representing the response of the  $j$ th unit selected equiprobably for the sample, we say:

$Y_j$  is a random variable whose distribution represents the possible values of the measured response variate for the  $j$ th unit in the sample of  $n$  units selected equiprobably from the respondent population, if the selecting and measuring processes were to be repeated over and over.

To pursue this idea of repeating over and over (or of *repetition*), we first discuss the set of all possible samples of size  $n$  that can be selected (without replacement) from a population of size  $N$ . There are  $N^{(n)}$  such samples if the *order* of selection is taken into account; for example, there are  $4^{(2)} = 12$  such samples of size 2 for a population of 4 elements.

**Example 6.1.1:** A respondent population of  $N = 4$  elements has the following integer *true*  $Y$ -values for its response variate: 1, 2, 4, 5 (so that:  $\bar{Y} = 3, S \approx 1.8257$ ).

Table 6.3.1 at the right lists the 12 possible ordered samples of size  $n = 2$  that can be selected; *equiprobable* selecting (EPS) is assumed. For discussion below, the values of the average ( $\bar{y}$ ) and sample error for each sample are also given.

Looking down the column of values for the *first* unit in the 12 samples, we see they are three copies of the values of the response variate of the elements of the respondent population; we see the same in the column of 12 values for the *second* unit. This result holds in general – in the set of  $N^{(n)}$  possible ordered samples of size  $n$ , selected from a population of size  $N$ , the  $j$ th unit comprises  $(N-1)^{(n-1)}$  copies of the values of the response variate of the elements of the respondent population. Hence, *under EPS*, the values of the unit in *any* position  $j$  of the set of all possible ordered samples of size  $n$  can be modelled like the values of the response variate of the elements of the respondent population – see equation (6.1.6) at the right; in statistics, it is more convenient to write equation the (6.1.6) as

$$\tau Y_j \sim N(\mu, \sigma_s), \quad j=1, 2, \dots, n, \quad \text{EPS} \quad \text{----(6.1.6)}$$

$$\text{response model (6.1.7)} \quad \tau Y_j = \mu + \tau R_j, \quad j=1, 2, \dots, n, \quad \tau R_j \sim N(0, \sigma_s), \quad \text{independent, EPS} \quad \text{----(6.1.7)}$$

**Table 6.1.1**

EPS of $n = 2$ units		
Sample	$\bar{y}$	Error
(1, 2)	1½	-1½
(2, 1)	1½	-1½
(1, 4)	2½	-½
(4, 1)	2½	-½
(1, 5)	3	0
(5, 1)	3	0
(2, 4)	3	0
(4, 2)	3	0
(2, 5)	3½	½
(5, 2)	3½	½
(4, 5)	4½	1½
(5, 4)	4½	1½

**NOTES:** 1. The context of Example 6.1.1 is *true* response variate values, indicated by the suffix T on the random variables  $\tau Y_j$  and  $\tau R_j$  (the **residuals**); conceptually, we are in the *first* sample ellipse in the schema at the bottom right overleaf on page 6.3 (see also Note 8 on page 6.6).

- Another difference between equations (6.1.1) and (6.1.6) is the subscript  $s$  on  $\sigma_s$ , which indicates the variation quantified by this (probabilistic) standard deviation arises from the *selecting* (or ‘sampling’) process.
- 2. Because any sample position contains *multiple* copies (not just one copy) of the values of the respondent population responses, its (data) standard deviation is slightly *smaller* than the respondent population (data) standard deviation  $S$ , which is represented by  $\sigma_s$  in the model. This discrepancy is of no consequence for the present discussion in the usual practical situation of  $N$  large and  $N \gg n$  (see also Note 4 on the facing page 6.5).

- 3. The second and third columns of Table 6.1.1 in Example 6.1.1 above illustrate two other matters.
  - Sample error – the difference between the value of the sample and respondent population attributes (here, averages) – has an *average* of *zero* over the set of all possible samples; stated another way, the average of the set of all possible sample averages is the respondent population average. This is what is meant by the statement, for example in the first sentence of the second paragraph of Figure 6.2, that, under EPS, the sample average is an *unbiased estimator* of the study population average. Symbolically, we write:

$$E(\tau \bar{Y}) = \bar{Y} \quad \text{or} \quad E(\tau \bar{Y}) - \bar{Y} = 0 \quad [\text{or, in the model: } E(\tau \bar{Y}) = \mu \quad \text{or} \quad E(\tau \bar{Y}) - \mu = 0]. \quad \text{----(6.1.8)}$$

This matter is illustrated more broadly in Appendix 3, starting on page 6.7 of this Figure 6.1.

– EPS is needed for unbiasedness when estimating  $\bar{Y}$  to ensure that all samples are *equally likely*, as reflected in calculating the *average* sample error.

- For the four  $Y$ -values of the response variate in this population, there *are* samples whose average is *equal* to the respondent population average,  $\bar{Y}$ ; however, there are many populations for which *no* sample of a given size (or, conceivably, of *any* size other than  $N$ ) has an average of  $\bar{Y}$ ; *i.e.*, *no* sample has *zero* sample error.
  - This is one reason a numerical Answer obtained by sampling is likely to be at least a little off the truth.

**Figure 6.1. FROM PROBABILITY TO STATISTICS: Modelling Sample Error... (continued 1)**

**4. A Probability Model for Selecting Equiprobably**

Example 6.1.1 on the facing page 6.4 illustrates properties of the set of all possible samples of a population under EPS; we now argue from this case to the process of selecting equiprobably (a sample) *over and over*. Provided ‘over and over’ means selecting a number of times that is *large* in relation to  $(N-1)^{(n-1)}$ , because any sample is *equally* likely, the set of samples will be (roughly) many copies of the set of all possible samples. Subject to this new caveat of not having *exactly* equal proportions of each possible sample in the set, together with a slightly smaller standard deviation (see Note 2 on the facing page 6.4), we can again use the model (6.1.7). Hence, for the behaviour of the sample *average* when selecting over and over, we use the model (6.1.3), rewritten above as equation (6.1.9).

$$\tau\bar{Y} \sim N(\mu, \sigma_s \sqrt{\frac{1}{n}}), \text{ EPS} \quad \text{-----(6.1.9)}$$

\* The idea of obtaining, when selecting over and over, roughly *equal* proportions of each member of the set of all possible samples is like tossing a fair coin over and over, where we expect to observe roughly *equal* proportions of heads and tails.

**NOTES:** 4. Note 2 and the paragraph above refer to the standard deviation of  $\tau Y_j$  as being slightly *smaller* than  $\sigma_s$ ; using a different approach in STAT 332, the standard deviation of  $\tau Y_j$  is actually as given in equation (6.1.10); the multiplier of  $(1-1/N)$  under a square root is obviously close to 1 in value for most populations encountered in practice, which have *many* elements.

$$s.d.(\tau Y_j) = \mathbf{S} \sqrt{1 - \frac{1}{N}} \quad \text{-----(6.1.10)}$$

$$s.d.(\tau\bar{Y}) = \mathbf{S} \sqrt{\frac{1}{n} - \frac{1}{N}} \quad \text{-----(6.1.11)}$$

● The consequence of equation (6.1.10) is that the standard deviation of  $\tau\bar{Y}$  is as given in equation (6.1.11) above; in the usual situation in practice, where  $N \gg n$ , this standard deviation becomes effectively that of equation (6.1.9), where  $\sigma_s$  is the model parameter representing the respondent population (data) standard deviation  $\mathbf{S}$ .

5. The assumption of a *normal* model for the shape of the histogram of the (true) response variate  $\mathbf{Y}$ -values in the respondent population can be relaxed if we can assume that the Central Limit Theorem will provide adequate *approximate* normality of the random variable  $\tau\bar{Y}$  representing the sample average  $\tau\bar{Y}$  under repetition.

**5. An EPS-Based Probability Model for Measuring**

The response model (6.1.7), with normal residuals, for equiprobable selecting over and over can be adapted to model the process of *measuring* (independently) over and over the value of a variate for a unit, for two reasons:

- \* We know from histograms like those for the paper thickness data in Figure 2.8b and the coin weights in Figure 2.9b in Part 2 that normal distributions are a reasonable model for values produced by (some) measuring processes.
- \* We can model measuring as a selecting process: possible values for the variate being measured could be written on slips of paper, the slips placed in a box and a slip selected by EPS; the number on the slip is then regarded as the ‘measured’ value; measuring over and over would be modelled as selecting slips equiprobably *with* replacement over and over.

We therefore write the response model for measuring as equation (6.1.12) at the right, where:

$${}_M Y = \tau + \delta + {}_M R, \quad {}_M R \sim N(0, \sigma_M), \text{ independent, EPS} \quad \text{-----(6.1.12)}$$

${}_M Y$  is a random variable whose distribution represents the possible values of the measurement of the response variate of a unit, if the measuring process were to be repeated over and over on this unit.

$\tau$  is a model parameter which represents the *true value* of the response variate of the unit being measured.

$\delta$  is a model parameter (called the *bias*) which represents the *inaccuracy* of the measuring process; the value of  $\delta$  *quantifies* the inaccuracy of the measuring process – as inaccuracy *increases* (*i.e.*, as accuracy *decreases*),  $\delta$  *increases*.

$R$  is a random variable (called the *residual*) whose distribution represents the possible *differences*, from the structural component of the model, of the value of the measurement of the response variate of the unit being measured, if the measuring process were to be repeated over and over on this unit.

$\sigma_M$  the *standard deviation* of the normal model for the distribution of the residual, is a model parameter (called the *variability*) which represents the *imprecision* of the measuring process and describes measuring variation if the measuring process were to be repeated over and over on a unit; the value of  $\sigma_M$  *quantifies* the imprecision of the measuring process – as imprecision *increases* (*i.e.*, as precision *decreases*),  $\sigma_M$  *increases*.

For a *calibrated* measuring system, where the inaccuracy has been (essentially) removed using a defined standard, the parameter  $\delta$  can be omitted and the model written as equation (6.1.13).

$${}_M Y = \tau + {}_M R, \quad {}_M R \sim N(0, \sigma_M), \text{ independent, EPS} \quad \text{-----(6.1.13)}$$

**NOTES:** 6. The absence of the subscript  $j$  on  ${}_M Y$  and  ${}_M R$  in the response models (6.1.12) and (6.1.13) is because we are modelling the process of measuring *once* the value of a variate for a unit.

7. The suffix  $M$  on  ${}_M Y$  (and  ${}_M R$ ) in the models (6.1.12) and (6.1.13) reminds us we are dealing with a *measured* value of a response variate, as distinct from  $\tau Y_j$  in equation (6.1.7), the corresponding *true* value for the  $j$ th unit in the sample.

**6. A Probability Model for Selecting Equiprobably and Measuring**

To obtain a response model for both selecting *and* measuring, two processes responsible for *incomplete* data and, hence, for *uncertain* Answers, we combine the response models (6.1.7) and (6.1.13); the resulting model [equation (6.1.16) below] describes the processes of EPS of  $n$  units from a respondent population (for which  $\mathbb{N} \gg n$ ) and measuring their response variate values once each with a *calibrated* measuring system.

Equation (6.1.7) is given again at the right, 
$${}_T Y_j = \mu + {}_T R_j, \quad j=1, 2, \dots, n, \quad {}_T R_j \sim N(0, \sigma_s), \quad \text{independent, EPS} \quad \text{----(6.1.7)}$$

and equation (6.1.13) is rewritten as equation (6.1.14), where  $\tau$  has been replaced 
$${}_M Y_j = {}_T Y_j + {}_M R_j, \quad j=1, 2, \dots, n, \quad {}_M R_j \sim N(0, \sigma_m), \quad \text{independent, EPS} \quad \text{----(6.1.14)}$$

by the *random variable*  ${}_T Y_j$  because we are measuring (once) the response of the  $j$ th unit in the sample *selected by EPS*;  ${}_M Y_j$  has a subscript  $j$  for the same reason. Also, the suffixes T and M on the residuals are to distinguish the residuals in the two models for selecting and measuring; these two sets of residuals become *one* set denoted  $R_j$  in the combined model (6.1.16) – see equation (6.1.17) below.

Then, writing  ${}_T Y_j$  in equation (6.1.14) as its expression in (6.1.7), we obtain 
$${}_M Y_j = \mu + {}_T R_j + {}_M R_j, \quad j=1, 2, \dots, n \quad \text{----(6.1.15)}$$

equation (6.1.15), which we rewrite as equation (6.1.16), where  $\sigma$  quantifies the *overall* 
$${}_M Y_j = \mu + R_j, \quad j=1, 2, \dots, n, \quad R_j \sim N(0, \sigma), \quad \text{independent, EPS} \quad \text{----(6.1.16)}$$

variation due to *both* selecting and measuring – see equation (6.1.17). 
$$R_j = {}_T R_j + {}_M R_j, \quad \sigma = \sqrt{\sigma_s^2 + \sigma_m^2} \quad \text{----(6.1.17)}$$

**NOTES:** 8. In this Figure 6.1, we have developed a model for selecting *and* measuring by combining separate models for the two processes; in doing so, we distinguish true and measured response variate values ( ${}_T Y_j$  and  ${}_M Y_j$ ). From now on in these Course Materials, the data we encounter are essentially *always* obtained by measuring and explicit consideration of the true value  ${}_T Y_j$  of equation (6.1.7) seldom occurs. For convenience, we will therefore usually write the model (6.1.16) as at the right *without* the suffix M.

$$Y_j = \mu + R_j, \quad j=1, 2, \dots, n, \quad R_j \sim N(0, \sigma), \quad \text{independent, EPS} \quad \text{----(6.1.16)}$$

- ${}_M Y_j$  and  $Y_j$  in equation (6.1.16) is the quantity defined early in Section 3 on the second side (page 6.4) of this Figure 6.1 – this definition involves *both* selecting and measuring; conceptually, we have reached the *second* sample ellipse in the schema at the lower right of the first side (page 6.3) of this Figure 6.1.

9. Based on the model (6.1.16), equation (6.1.18) at the right describes the behaviour of the random variable representing the sample average,  $\bar{Y}$ , subject to variation from *both* equiprobable selecting and measuring.

$$\bar{Y} \sim N(\mu, \sigma \sqrt{\frac{1}{n}}), \quad \text{EPS} \quad \text{----(6.1.18)}$$

- Despite the apparent similarity of equations (6.1.3) and (6.1.18), their derivations and the interpretation of  $\bar{Y}$ ,  $\mu$  and  $\sigma$  differ in statistically vital ways.

**7. Appendix 1: A Comparison of Approaches to Modelling the Behaviour of  $\bar{Y}$**

Three approaches to modelling the behaviour of the sample average, as a basis for estimating the respondent population average, are discussed in this Figure 6.1; this Appendix 1 presents an overview of the strengths and weaknesses of each.

- \* *Probability modelling* starts with the normal model (6.1.1) and then uses probability theory for random variables to obtain the model (6.1.3); the weakness of this approach for *statistics* is that the ideas of Sections 2 to 6, and the assumptions they make *explicit*, are only *implicit*. It is likely to be unclear to the beginning student why statements like equations (6.1.1) and (6.1.3) can be used to describe the real-world processes of selecting (and measuring, if it is even considered in this approach) and it is easy for statistical practice based on the probability approach to overlook essential components of the Plan.

$$Y \sim N(\mu, \sigma) \quad \text{----(6.1.1)}$$

$$\bar{Y} \sim N(\mu, \sigma \sqrt{\frac{1}{n}}) \quad \text{----(6.1.3)}$$

- \* The *response model* (6.1.16), as developed on the basis of the *same* probability models (6.1.1) and (6.1.3), recognizes explicitly Question formulation (via consideration of the target, study and respondent populations) and properties of the selecting and measuring processes needed as a basis for the relevant probability models, written as the response models (6.1.7) and (6.1.14). This approach therefore overcomes the weaknesses of the probability modelling approach, albeit at the cost of more detail in presentation; the final probability expression for the behaviour of  $\bar{Y}$  involving variation from *both* selecting and measuring, is equation (6.1.18). A weakness of this approach is that it does not recognize explicitly the finite size ( $\mathbb{N}$  elements) of the respondent population.

$${}_T Y_j = \mu + {}_T R_j, \quad j=1, 2, \dots, n, \quad {}_T R_j \sim N(0, \sigma_s), \quad \text{indep., EPS} \quad \text{----(6.1.7)}$$

$${}_M Y_j = {}_T Y_j + {}_M R_j, \quad j=1, 2, \dots, n, \quad {}_M R_j \sim N(0, \sigma_m), \quad \text{indep., EPS} \quad \text{----(6.1.14)}$$

$$Y_j = \mu + R_j, \quad j=1, 2, \dots, n, \quad R_j \sim N(0, \sigma), \quad \text{indep., EPS} \quad \text{----(6.1.16)}$$

$$\bar{Y} \sim N(\mu, \sigma \sqrt{\frac{1}{n}}), \quad \text{EPS} \quad \text{----(6.1.18)}$$

- \* The so-called *design-based* approach developed in STAT 332, based on the equiprobable selecting of the sampling protocol, incorporates the finite size ( $\mathbb{N}$  elements) of the respondent population, but it cannot easily take direct account of measurement error. Its expression for the behaviour of  ${}_T \bar{Y}$  is equation (6.1.19); the *approximate* normality comes from the Central Limit Theorem approximation.

$${}_T \bar{Y} \doteq N(\bar{\mathbf{Y}}, \mathbf{S} \sqrt{\frac{1}{n} - \frac{1}{\mathbb{N}}}), \quad \text{EPS} \quad \text{----(6.1.19)}$$

Estimating  $\bar{Y}$  using the models (6.1.3), (6.1.18) and (6.1.19) requires an *estimate* of  $\sigma$  or  $\mathbf{S}$ , customarily taken as the *sample* stan-

**Figure 6.1. FROM PROBABILITY TO STATISTICS: Modelling Sample Error... (continued 2)**

standard deviation,  $s$ . Because  $s$  is calculated from *measured* values of the response variate, its value *includes* variation due to measuring *as well as* selecting; this inclusion is modelled in (6.1.18) but is a *fortuitous* benefit for (6.1.3) and (6.1.19).

Thus, the model (6.1.18), derived from the response model (6.1.16), provides a basis for statistical practice that is missing from (6.1.3); the additional insight from (6.1.19) is that the multiplier of  $\mathbf{S}$  in the standard deviation of  $\tau\bar{Y}$  has weak dependence on  $\mathbf{N}$ ; ignoring  $\mathbf{N}$  in (6.1.18) is justified provided the sample size is a *small proportion* of the respondent population size ( $n \ll \mathbf{N}$ ).

**8. Appendix 2: Using the Response Model for Measuring**

When investigating the inaccuracy and imprecision of a measuring system [often as part of broader data-based investigating to answer Question(s) of interest], a common approach is to make  $m$  independent measurements of the same value. The appropriate response models are shown

at the right as equations (6.1.20) and  ${}_mY_j = \tau + \delta + {}_mR_j, \quad j = 1, 2, \dots, m, \quad {}_mR_j \sim N(0, \sigma_M), \quad \text{independent, EPS} \quad \text{----(6.1.20)}$

(6.1.21); the former is applicable when  ${}_mY_j - \tau = \delta + {}_mR_j, \quad j = 1, 2, \dots, m, \quad {}_mR_j \sim N(0, \sigma_M), \quad \text{independent, EPS} \quad \text{----(6.1.20)}$

measuring a *known* value to investigate  ${}_mY_j = \tau + {}_mR_j, \quad j = 1, 2, \dots, m, \quad {}_mR_j \sim N(0, \sigma_M), \quad \text{independent, EPS} \quad \text{----(6.1.21)}$

both inaccuracy *and* imprecision, whereas the latter is used when investigating only imprecision; the two forms of (6.1.20) differ in whether the measured value or its difference from the true value is taken as the response variate – one or the other may be more convenient, depending on context.

**Example 6.1.2:** A medical laboratory purchases a standard, certified to contain 200 mg/dL of cholesterol, to calibrate its process for measuring serum cholesterol levels; data from 15 measurements taken over an 8-hour day were:

Table	Measured value ( ${}_mY_j$ )	194	196	202	198	195	188	203	200	196	195	198	202	203	201	199	Av.	S.d.
<b>6.1.2:</b>	Measured – true value ( ${}_mY_j - \tau$ )	-6	-4	2	-2	-5	-12	3	0	-4	-5	-2	2	3	1	-1	-2	4.1231

These data show that: the estimate of measuring *inaccuracy* (represented in the model by  $\delta$ ) is  $-2$  mg/dL; the estimate of measuring *imprecision* (represented in the model by  $\sigma_M$ ) is 4.1 mg/dL.

With inaccuracy estimated to be about half the magnitude of imprecision, the former may have little practical importance here, but more detailed interpretation of these data requires extra-statistical knowledge.

**NOTE:** 10. Under the idealizations of the models (6.1.20) and (6.1.21), it can be shown that, for measuring processes:

- inaccuracy is *unaffected* by averaging – looking ahead, one average is the estimate of the *intercept* of the centred form of the straight-line model in simple linear regression (see page 13.19 in STAT 221 Figure 13.3);
- inaccuracy is *removed* by differencing – differences occur when comparing and when calculating (data) standard deviations and the *slope* of the straight-line regression model (see page 13.19 in STAT 221 Figure 13.3).

**9. Appendix 3: The Set of All Possible Samples from a Population**

Section 3 introduced the idea of all possible samples that can be selected from a population. The sets of ten and nine histograms on pages 6.8 and 6.9 illustrate this idea for a population of  $\mathbf{N} = 10$  elements with response variate values  $\mathbf{Y} = 1, 2, \dots, 10$ . These histograms, for sample averages and sample (data) standard deviations, show important properties of the selecting process.

- \* The sample average takes a value which is determined by the particular sample selected; these values, over the set of all possible samples of a given size, form a *distribution*; under EPS, *this is the distribution of the random variable*  $\bar{Y}$ , a result that could be called **The Fundamental Theorem of Statistics** by analogy with The Fundamental Theorem of Calculus.
  - The mean of this distribution is the *population average* – this is the *unbiasedness* under EPS of the sample average as an estimator of the population average (see the second column of Table 6.1.3 in Appendix 4 on page 6.10).
  - The standard deviation of this distribution is given by equation (6.1.11) except that, because  $s.d.(\bar{Y})$  is a *probabilistic* standard deviation and  $s_{\bar{Y}}$  (given in the third column of Table 6.1.3 in Appendix 4) is a *data* standard deviation, the former is  $\sqrt{k/(k-1)}$  times the latter, where  $k$  is  $\binom{\mathbf{N}}{n}$ , the number of possible (unordered) samples of size  $n$ .
- \* The distribution of sample averages reminds us that for the *particular* sample selected when executing the Plan for an investigation, its sample error is *unknown* – its attribute may be close to the population attribute or (sometimes) it may not be; that is, an Answer has *sampling uncertainty*, meaning we cannot *know* how close an Answer obtained by sampling is to the truth.
- \* The *decreasing variation* of the values of sample averages with *increasing* sample size, visible as the decreasing width down the page of the histograms of sample averages, illustrates *sampling imprecision* decreasing with increasing *sample size*; *i.e.*, under EPS, increasing sample size decreases the average magnitude of sample error and so decreases *sampling uncertainty*.
- \* The change in distribution *shape*, most noticeable as the sample size increases from 1 to 2, illustrates the *idea* behind the CLT – a *non-normal* (uniform) distribution starting to show the central ‘heaping’ of the normal distribution.
  - The heaping, which first appears when  $n = 2$ , persists up to  $n = 8$  and then disappears; this is likely the effect of *dependence* among samples, which consist of most of the in the population, becoming dominant over the heaping process reminiscent of the Central Limit Theorem.

(continued overleaf)

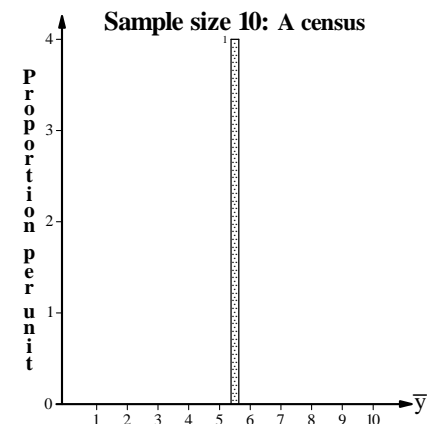
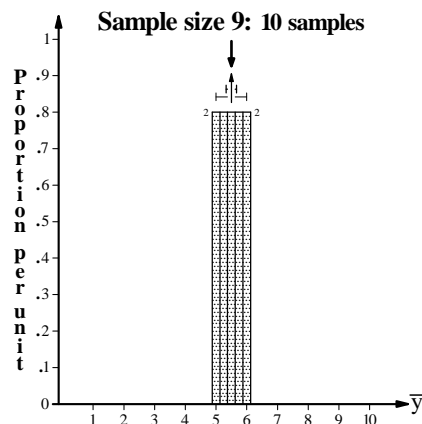
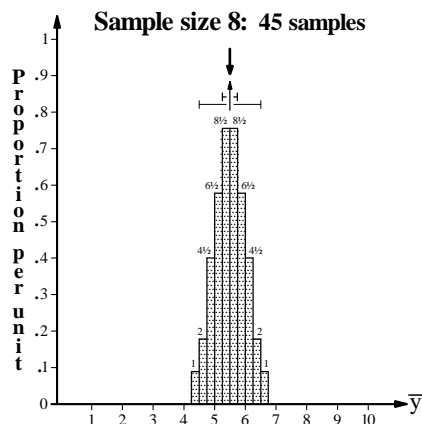
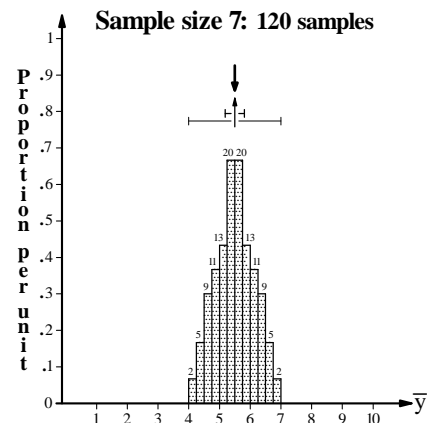
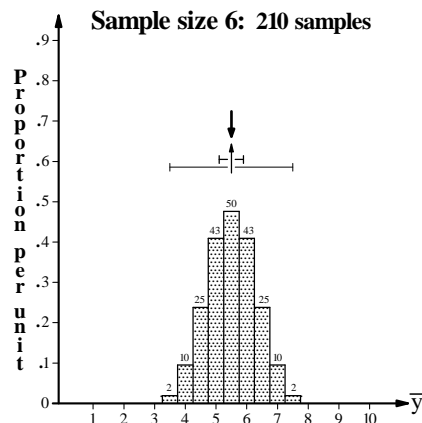
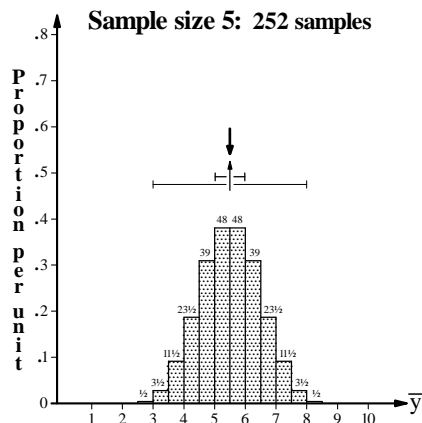
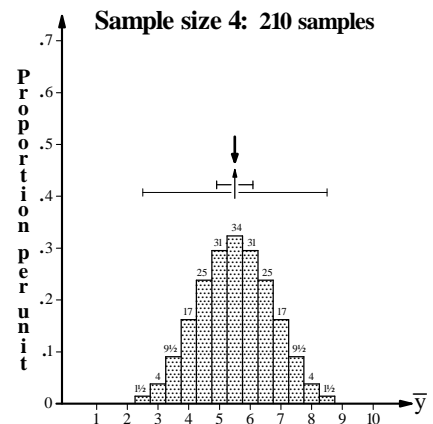
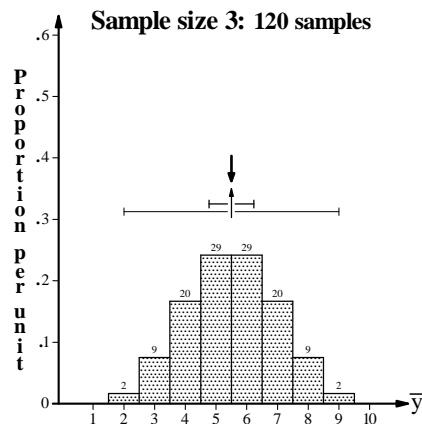
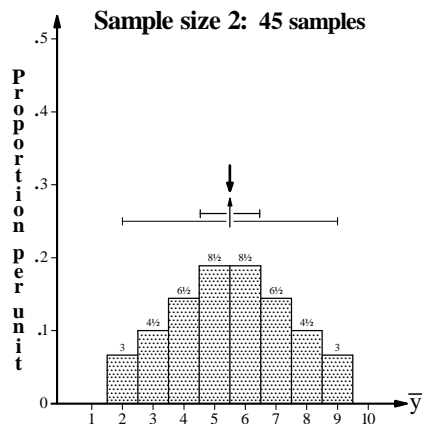
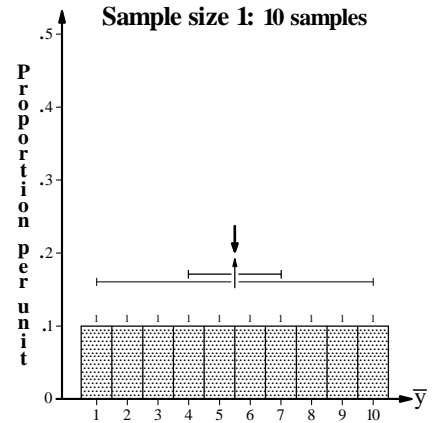
**9. Appendix 3: The Set of All Possible Samples from a Population (continued)**

These ten histograms show the distribution of the *averages* of all possible samples of a given size that can be selected without replacement from a population of  $N = 10$  elements, whose response variate values are  $Y = 1, 2, \dots, 10$ . [A histogram of these ten population *responses* would be like that given at the right for the averages of the ten possible samples of size 1.]

The decreasing vertical axis scale unit down the page reduces the visual impact of increasing histogram height; bar *frequencies* are given by the numbers at the tops of the bars.

Above each histogram, except the one at the bottom right of the page:

- the upper bold arrow ( $\downarrow$ ) indicates the value of the *population average*  $\bar{Y} = 5.5$ ;
- the lower arrow ( $\uparrow$ ) indicates the *average* of the set of *sample averages* ( $\bar{\bar{y}} = 5.5$ );
- the upper bar ( $\text{---}$ ) crossing the lower arrow has a length equal to the value of the (data) *standard deviation* of the set of sample averages;
- the lower such bar shows the *range* of the set of sample averages.



**Figure 6.1. FROM PROBABILITY TO STATISTICS: Modelling Sample Error... (continued 3)**

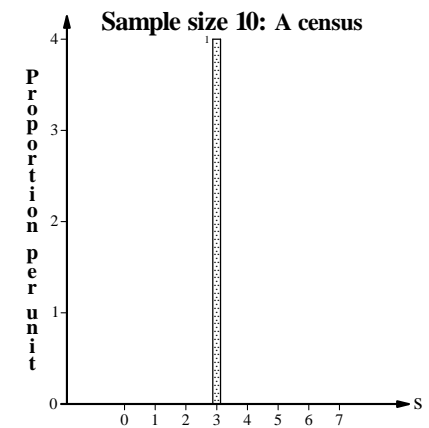
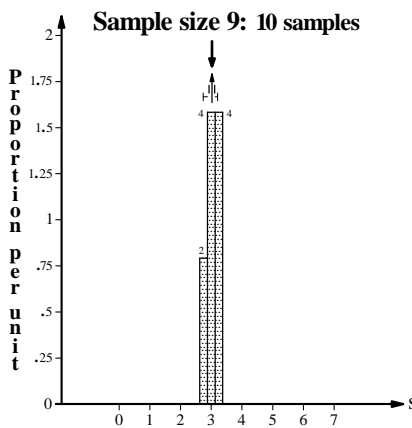
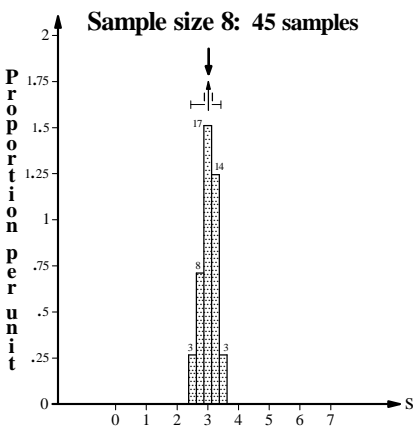
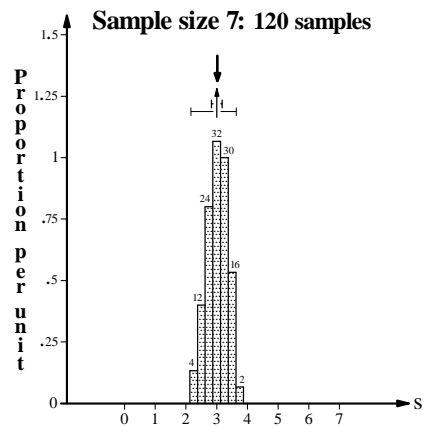
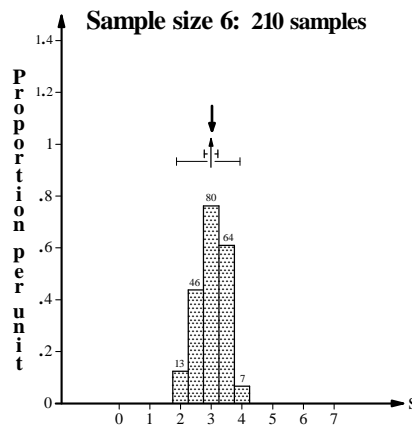
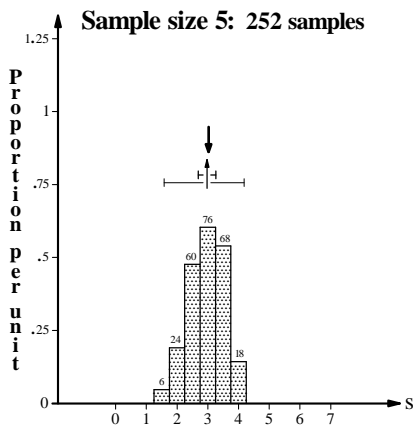
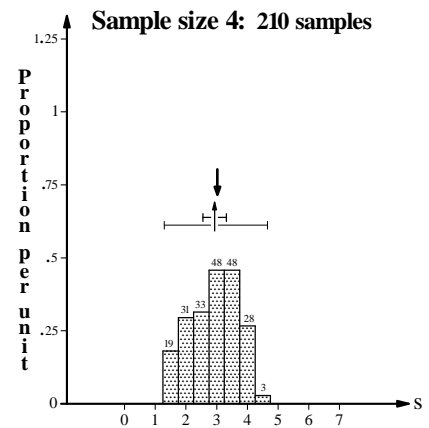
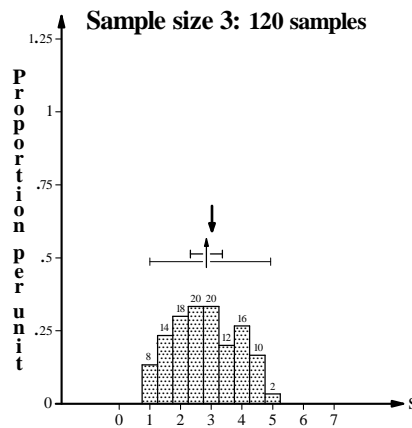
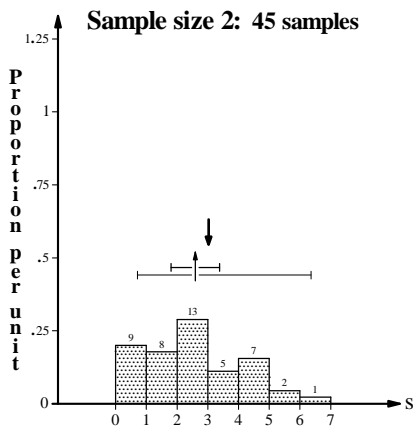
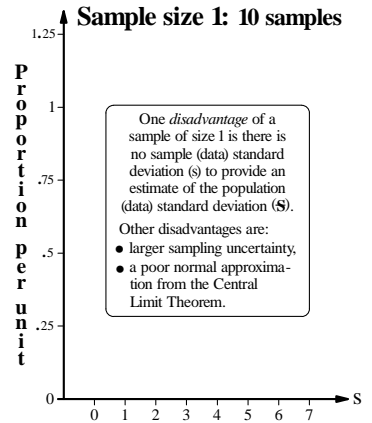
These nine histograms show the distribution of the (data) *standard deviations* of all possible samples of a given size that can be selected without replacement from a population of  $N = 10$  elements, whose response variate values are  $Y = 1, 2, \dots, 10$ .

[A histogram of these ten population *responses* would be like that given at the top right of the facing page 5.92 for the averages of the ten possible samples of size 1.]

The decreasing vertical axis scale unit down the page reduces the visual impact of increasing histogram height; bar *frequencies* are given by the numbers at the tops of the bars.

Above each histogram, except the one at the bottom right of the page:

- the upper bold arrow ( $\downarrow$ ) indicates the value of the *population* standard deviation  $S = 3.027\ 650$ ;
- the lower arrow ( $\uparrow$ ) indicates the *average* ( $\bar{s}$ ) of the set of *sample* (data) standard deviations;
- the upper bar ( $\text{---}$ ) crossing the lower arrow has a length equal to the value of the (data) *standard deviation* of the set of sample (data) standard deviations;
- the lower such bar shows the *range* of the set of sample (data) standard deviations.









**11. Appendix 5: Sampling Protocols Beyond EPS**

Equiprobable (or simple random) selecting of units consisting of *individual elements* from an *unstratified* (respondent) population is useful for modelling the selecting process but, in practice, more complex sampling protocols are used. Two such protocols are:

- \* **cluster selecting:** selecting equiprobably units from the (respondent) population that are *groups* of elements – clusters may be of *equal* size (e.g., cardboard boxes of 24 cans of soup) or *unequal* size (e.g., households);
- \* **stratified selecting:** subdividing the (respondent) population into groups (called **strata**) so that elements *within* a stratum have *similar* response variate values ('homogeneity of strata') and elements in different strata *differ* as much as practicable from each other; the sample is obtained by equiprobable selecting of units consisting of individual elements from *each* stratum.

Example 6.1.3 below illustrates the effects of *clustering* and *stratifying* on *precision*, also bearing in mind that:

- *clustering* is commonly used because of the availability of a clustered *frame*, thereby avoiding the cost of generating a frame as part of the investigating – a **frame** can be thought of as a *list* of the (respondent) population sampling units (here, clusters);
- *stratifying* is commonly used because it provides (the often useful additional) *subdivision* of Answers by stratum (e.g., in Canada, the breakdown by province and territory of the national unemployment rate).

**Example 6.1.3:** A respondent population of  $N = 4$  elements (or units) has the following integer  $\mathbf{Y}$ -values for its response variate:

1, 2, 4, 5 [so that the population average and (data) standard deviation are:  $\bar{Y} = 3$ ,  $S = 1.8257$ ];

we examine the *sampling imprecision*, under equiprobable selecting (EPS) with a sample size of  $n = 2$ , of the random variable  $\bar{Y}$  whose values are the sample average  $\bar{y}$ , as an estimator of  $\bar{Y}$  using three sampling protocols:

- EPS of two units, each consisting of one element, from the *unstratified* population;
- EPS of one *cluster*, of size  $L = 2$  elements, from the *unstratified* population;
- EPS of one unit, consisting of one element, from each of two *strata* of size  $N_1 = N_2 = 2$ .

Note that *each* estimator is *unbiased*, because  $E(\bar{Y}) = \bar{Y}$  or  $E(\bar{Y}) - \bar{Y} = 0$  [the *average* sample error is *zero*].

**Table 6.1.5**

**Unstratified population  
EPS of two elements**

Sample	$\bar{y}$	Error
(1, 2)	1½	-1½ large
(1, 4)	2½	-½ medium
(1, 5)	3	0 small
(2, 4)	3	0 small
(2, 5)	3½	½ medium
(4, 5)	4½	1½ large

Designation of sample error as *large*, *medium* or *small* is *only* for convenience in the context of Example 6.1.3.

**Table 6.1.6a**

**Unstratified population  
Clusters: [1, 2], [4, 5]  
EPS of one cluster (L = 2)**

Sample	$\bar{y}$	Error
(1, 2)	1½	-1½ large
(4, 5)	4½	1½ large

**Table 6.1.7a**

**Stratified population  
Strata: [1, 2], [4, 5]  
EPS of one element per stratum**

Sample	$\bar{y}$	Error
(1, 4)	2½	-½ medium
(1, 5)	3	0 small
(2, 4)	3	0 small
(2, 5)	3½	½ medium

**Table 6.1.6b**

**Unstratified population  
Clusters: [1, 4], [2, 5]  
EPS of one cluster (L = 2)**

Sample	$\bar{y}$	Error
(1, 4)	2½	-½ medium
(2, 5)	3½	½ medium

**Table 6.8.7b**

**Stratified population  
Strata: [1, 4], [2, 5]  
EPS of one element per stratum**

Sample	$\bar{y}$	Error
(1, 2)	1½	-1½ large
(1, 5)	3	0 small
(2, 4)	3	0 small
(4, 5)	4½	1½ large

Example 6.1.3 illustrates for us that:

- The effect on (sampling) imprecision of clustering and of stratifying depends on how each is *implemented* in the Plan – that is, it depends on this component of the *sampling protocol*.
- Clustering and stratifying affect precision by determining which of the *possible* samples of size  $n$  have non-zero selecting probabilities.
- Decreased imprecision is favoured by *heterogeneity of clusters* but by *homogeneity of strata* with respect to the response(s) of interest.
  - In the middle and right-hand columns of Example 5.8.3, heterogeneity increases *down* the three clustered sampling protocols, homogeneity increases *up* the three stratified protocols.

**Table 6.1.6c**

**Unstratified population  
Clusters: [1, 5], [2, 4]  
EPS of one cluster (L = 2)**

Sample	$\bar{y}$	Error
(1, 5)	3	0 small
(2, 4)	3	0 small

**Table 6.1.7c**

**Stratified population  
Strata: [1, 5], [2, 4]  
EPS of one element per stratum**

Sample	$\bar{y}$	Error
(1, 2)	1½	-1½ large
(1, 4)	2½	-½ medium
(2, 5)	3½	½ medium
(4, 5)	4½	1½ large

[As an exercise, quantify the sample error *variation* by calculating the relevant (data) *standard deviation* for each of the seven sampling protocols; comment on what is illustrated by the values obtained.]

- There is a sense in which clustering is *passively* accepted in the interests of reducing investigation cost, whereas stratifying is *actively* imposed by the investigator(s) on [or may be a natural feature of] the study (or respondent) population.
- While EPS from an *unstratified* population implies *equal* inclusion probabilities for all population *units*, the converse does *not* hold – in the three clustered and three stratified sampling protocols, all *elements* have equal inclusion probabilities but all six *samples* of size 2 are *not* equally probable [four samples and two samples (respectively) have *zero* selecting probability].