

Figure 2.10. UNSTRATIFIED POPULATIONS: One-Stage EPSWOR of Individual Elements
Estimating a Proportion or a Frequency

The theory developed in Figure 2.3 is concerned with equiprobable selecting (EPS) when the response, \mathbf{Y} , is the *measured* value of some *quantitative* respondent population response variate, and we wish to estimate a respondent population attribute such as an average or total. In this Figure 2.10, we develop the corresponding theory for the *proportion* of respondent population elements with some *qualitative* characteristic. Examples of proportions of practical interest are:

- * the proportion (or percentage) of voters who support a given political party (for a political poll);
- * the proportion of TV viewers watching a particular program (for TV ratings);
- * the proportion of sales of an item of a specified brand (for market research);
- * the proportion of trees in a forest with trunks greater than 30 cm in diameter (for a lumber company).

With appropriate changes in notation, the theory of EPS for a proportion can be derived from the theory for an average given in Figure 2.3. We deal here with a response in *two* categories [C and cC (the complement of C), say] for the characteristic of the respondent population elements; in Figure 2.12, we extend the ideas to a response in *more than two* categories.

- 1. Notation.** Define the indicator variate \mathbf{V} such that: $\mathbf{V}_i = 1$ if respondent population element i has the characteristic C ,
 $\mathbf{V}_i = 0$ if respondent population element i has the characteristic cC .

Table 2.10.1:	QUANTITY.....	RESPONDENT POPULATION	SAMPLE [MODEL].....
	Size (elements/units)	\mathbf{N}	n
	Indicator variate	\mathbf{V}_i ($i = 1, 2, \dots, \mathbf{N}$)	v_j ($j = 1, 2, \dots, n$) [r.v.s are V_j with value v_j]
	Proportion	$\mathbf{P} = \frac{1}{\mathbf{N}} \sum_{i=1}^{\mathbf{N}} \mathbf{V}_i = \bar{\mathbf{V}}$	$p = \frac{1}{n} \sum_{j=1}^n v_j = \bar{v}$ [r.v. is P with value p ($= \bar{v}$)]
	Number (or frequency)	\mathbf{NP}	(number in sample $= np$)

The response of element j obtained by EPS from the population of \mathbf{N} elements can be modelled by the random variable V_j with probability function and mean as in equation (2.10.1) at the right. It follows that the random variable P is an unbiased estimator of \mathbf{P} – see equation (2.10.2), the analogue of equation (2.3.8) on page 2.41 of Figure 2.3.

$$\begin{array}{c|cc} v_j & 1 & 0 \\ \hline f(v_j) & \mathbf{P} & 1-\mathbf{P} \end{array}; \quad E(V_j) = 1 \times \mathbf{P} + 0 \times (1-\mathbf{P}) = \mathbf{P} \quad \text{-----(2.10.1)}$$

$$E(P) = E\left(\frac{1}{n} \sum_{j=1}^n V_j\right) = \frac{1}{n} \sum_{j=1}^n E(V_j) = \frac{1}{n} n E(V_j) = \frac{1}{n} n \mathbf{P} = \mathbf{P} \quad \text{-----(2.10.2)}$$

2. Estimating \mathbf{P} , the Respondent Population Proportion

We want both a value (or *point estimate*) for this population attribute *and* a measure of the uncertainty of the estimate, for which we use a confidence interval in STAT 332.

Using the indicator variate \mathbf{V} , we have shown above that a proportion becomes a special case of an average. Because it was established in equation (2.3.8) at the top of page 2.41 in Figure 2.3 that the random variable \bar{Y} representing the sample *average* under EPS is an unbiased estimator of $\bar{\mathbf{Y}}$, the respondent population average, the random variable P representing the sample *proportion* under EPS is an unbiased estimator of \mathbf{P} , the respondent population proportion [as in equation (2.10.2) above].

To find the standard deviation of \mathbf{P} , we note first that because the population of \mathbf{V} s consists only of 1s and 0s, both population sums $\sum \mathbf{V}_i$ and $\sum \mathbf{V}_i^2$ are equal to \mathbf{NP} , the number of elements with characteristic C in the respondent population; the corresponding *sample* sums are both equal to np ($= np$); we also use the notation $\mathbf{Q} = 1 - \mathbf{P}$.

$$\text{Hence: } \mathbf{S} = \sqrt{\frac{1}{\mathbf{N}-1} \sum_{i=1}^{\mathbf{N}} (\mathbf{V}_i - \bar{\mathbf{V}})^2} = \sqrt{\frac{1}{\mathbf{N}-1} [\sum_{i=1}^{\mathbf{N}} \mathbf{V}_i^2 - \mathbf{N} \bar{\mathbf{V}}^2]} = \sqrt{\frac{1}{\mathbf{N}-1} [\mathbf{NP} - \mathbf{NP}^2]} = \sqrt{\frac{\mathbf{N}}{\mathbf{N}-1} \mathbf{P}(1-\mathbf{P})} = \sqrt{\frac{\mathbf{N}}{\mathbf{N}-1} \mathbf{PQ}}, \quad \text{-----(2.10.3)}$$

$$\text{so that: } s.d.(P) = \mathbf{S} \sqrt{\frac{1}{n} - \frac{1}{\mathbf{N}}} = \sqrt{\frac{\mathbf{N}}{\mathbf{N}-1} \mathbf{PQ} \left(\frac{1}{n} - \frac{1}{\mathbf{N}}\right)} \quad (= \sqrt{\frac{\mathbf{N}-n}{\mathbf{N}-1} \mathbf{PQ} \frac{1}{n}} = \sqrt{\frac{\mathbf{N}-n}{\mathbf{N}} \frac{\mathbf{N}}{\mathbf{N}-1} \mathbf{PQ} \frac{1}{n}}) \quad (= \text{the s.d. of } P). \quad \text{-----(2.10.4)}$$

$$\text{Also: } s = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (v_j - \bar{v})^2} = \sqrt{\frac{1}{n-1} [\sum_{j=1}^n v_j^2 - n \bar{v}^2]} = \sqrt{\frac{1}{n-1} [np - np^2]} = \sqrt{\frac{n}{n-1} p(1-p)} = \sqrt{\frac{n}{n-1} pq} \quad (q = 1-p), \quad \text{-----(2.10.5)}$$

$$\text{so that: } \hat{s}.d.(P) = s \sqrt{\frac{1}{n} - \frac{1}{\mathbf{N}}} = \sqrt{\frac{n}{n-1} pq \left(\frac{1}{n} - \frac{1}{\mathbf{N}}\right)} \quad (= \sqrt{\frac{\mathbf{N}-n}{\mathbf{N}} pq \frac{1}{n-1}} = \sqrt{\frac{\mathbf{N}-n}{\mathbf{N}} \frac{n}{n-1} pq \frac{1}{n}}) \quad (= \text{the estimated s.d. of } P). \quad \text{-----(2.10.6)}$$

$$\text{If } \mathbf{N} \gg n, \text{ equation (2.10.6) yields the approximate equation (2.10.7).} \quad \hat{s}.d.(P) \approx \sqrt{\frac{pq}{n-1}} \quad \text{-----(2.10.7)}$$

$$\text{The coefficient of variation (c.v.) of } P \text{ [a measure of relative precision] is:} \quad c.v.(P) \equiv \frac{s.d.(P)}{\mathbf{P}} = \sqrt{\frac{\mathbf{N}}{\mathbf{N}-1} \left(\frac{1}{\mathbf{P}} - 1\right) \left(\frac{1}{n} - \frac{1}{\mathbf{N}}\right)}. \quad \text{-----(2.10.8)}$$

$$\text{If } \mathbf{N} \gg n, \text{ equation (2.10.8) yields the approximate equation (2.10.9).} \quad c.v.(P) \approx \sqrt{\frac{1}{n} \left(\frac{1}{\mathbf{P}} - 1\right)}; \quad \hat{c}.v.(P) \approx \sqrt{\frac{1}{n} \left(\frac{1}{p} - 1\right)} \quad \text{-----(2.10.9)}$$

- NOTES:**
- The rightmost two expressions in brackets () in equation (2.10.4) and in equation (2.10.6) have the symbols rearranged only to illustrate discussion in Notes 7 and 8 below and on page 2.83; the *working* expressions in these two equations are their *second* square roots on their right-hand sides.
 - Equation (2.10.4) shows that the magnitude of the standard deviation of P is determined by the values two respondent attributes, \mathbf{PQ} and \mathbf{N} , and n , the sample size; because typically $\mathbf{N} \gg n$, n usually dominates \mathbf{N} in *their* effect on the standard deviation of P
 - The weak (often negligible) effect of \mathbf{N} on precision of sample estimates (under EPS) is contrary to many people's intuition; *i.e.*, it seems counterintuitive that *population size* (like 30 million Canadians or 300 million Americans) usually makes little or no contribution to the limitation imposed on an Answer by sample error.
 - $\mathbf{PQ} = \mathbf{P}(1 - \mathbf{P})$ is effectively, of course, just the attribute \mathbf{P} .
 - The expression for $s.d.(P)$ in equation (2.10.4) overleaf on page 2.81 is potentially useful in three ways:
 - to assess the precision of the estimator P [if \mathbf{N} , n and \mathbf{P} are known];
 - to calculate the sample size needed to attain a specified precision [if \mathbf{N} and \mathbf{P} are known – see Figure 2.12];
 - to compare the precision of P with that of *other* estimators of \mathbf{P} .
 - When $n = \mathbf{N}$, $s.d.(P) = 0$, reminding us that, at least in principle, sample error is eliminated in a census of the respondent population.
 - Of course, study error, non-response error and measurement error must still be adequately managed in a census.
 - The random variable S^2 [with value s^2 taken as (the value of) the *square* of the sample (data) standard deviation under EPS] is an *unbiased* estimator of \mathbf{S}^2 , the *square* of the respondent population (data) standard deviation \mathbf{S} [as shown near the bottom of page 2.41 in Figure 2.3], so it follows that $n\mathbf{PQ}/(n-1)$ is an unbiased estimator of $\mathbf{NPQ}/(\mathbf{N}-1)$.
 - This unbiasedness, a staple of survey sampling theory, is of interest statistically but is of limited relevance in the context of this Figure 2.10, because the unbiasedness is lost as a consequence of the *square roots* in the rightmost terms in equations (2.10.5) and (2.10.3) – see also Appendix 5 on pages 2.46 and 2.47 in Figure 2.3.
 - Equations (2.10.8) and (2.10.9) overleaf at the bottom of page 2.81 show that the *relative* precision *decreases* [*i.e.*, $s.d.(P)$ becomes *larger* relative to \mathbf{P}] as \mathbf{P} decreases, reminding us that (very) large sample sizes are needed to estimate small proportions precisely.
 - In real-world sample surveys to estimate a proportion or frequency, the sample is selected probabilistically from a (finite respondent) population. This raises the possibility of two *probability models*, which are often presented as an urn containing balls of two colours, red and black say. [Statistically, the urn is the population with the balls as its elements, the two colours being the response variate 'values' C and C .]
 - For EPS *without* replacement, the model is *hypergeometric*, whose probability function for the *number* of red balls in n selections (*i.e.*, in a sample of size n) is equation (2.10.10) at the right; \mathbf{N} is the number of balls in the urn, \mathbf{R} are red (and $\mathbf{N} - \mathbf{R}$ are black). Unsurprisingly, the mean and (probabilistic) standard deviation of the hypergeometric distribution are, respectively, $n\mathbf{P}$ and equation (2.10.4) overleaf on page 2.81 with an additional n outside the square root to deal with a frequency rather than a proportion; our \mathbf{P} is (of course) \mathbf{R}/\mathbf{N} and the terms multiplying \mathbf{PQ} under the square root may be written $(\mathbf{N} - n)/(\mathbf{N} - 1)$ as in the third expression in equation (2.10.4).
 - It is sometimes mathematically convenient to assume selecting is *with* replacement, with the appropriate probability model then being *binomial* instead of hypergeometric – selecting balls with replacement from the urn is an instance of Bernoulli trials; the probability function for the *number* of red balls in n selections (*i.e.*, in a sample of size n) is equation (2.10.11) at the right – the binomial parameter (a *Greek* letter) π is the probability of 'success,' here, selecting a red ball so $\pi = \mathbf{R}/\mathbf{N}$ (our \mathbf{P}).
 The mean and standard deviation of the binomial distribution are $n\pi$ and $\sqrt{n\pi(1-\pi)}$ for the *number* of red balls in n selections, and π and $\sqrt{\pi(1-\pi)/n}$ for the *proportion* of red balls in n selections. Comparing this second standard deviation with the third expression in equation (2.10.4) on page 2.81, we see that selecting *without* replacement introduces the additional term $(\mathbf{N} - n)/(\mathbf{N} - 1)$ – mathematically, we can think of this term as dealing with the (slight) *lack of probabilistic independence* when selecting *without* replacement compared to selecting *with* replacement.
- Other comparisons are pursued in Note 8 below and on the facing page 2.83.
- The order of topics in most current introductory statistics texts means that, before students can encounter ideas like $s.d.(Y)$ in equation (2.3.9) near the top of page 2.41 in Figure 2.3 or $s.d.(P)$ in equation (2.10.4) on page 2.81, they are presented with two (likely new) ideas.

(continued)

Figure 2.10. UNSTRATIFIED POPULATIONS: One-Stage EPSWOR of Individual Elements Estimating a Proportion or a Frequency (continued 1)

NOTES: 8. ● For (much beloved of statisticians) n probabilistically independent identically distributed (abbreviated iid) random variables (Y_i , say) each with mean μ and standard deviation σ , the random variable \bar{Y} that is their *average* has the *same* mean μ but a (smaller) standard deviation of σ/\sqrt{n} [equation (2.10.12)].

$$s.d.(\bar{Y}) = \sigma/\sqrt{n} \equiv \sigma\sqrt{\frac{1}{n}} \quad \text{----(2.10.12)}$$

- The standard deviation for a binomial *proportion* also has \sqrt{n} in its denominator; as discussed on the facing page 2.82 in the second bullet (●) of Note 7.

However, with the accolades it tends to receive from statisticians, this useful insight that averages are *less* variable than their components by a factor of \sqrt{n} can become an ideal that other relevant statistical theory should try to emulate; in the case of survey sampling theory, pursuit of this ‘ideal’ is detrimental to statistics and its teaching in three ways.

- * Equations (2.3.9) and (2.3.17) near the top and bottom of page 2.41 in Figure 2.3 are written in their second (equivalent) forms shown as equations (2.10.13) and (2.10.14) at the right, as are equations (2.10.4) and (2.10.6) from page 2.81, where $f = n/N$; equation (2.10.6) is even rewritten as (the inelegant) equation (2.10.15) to yield (the mandatory?) $1/n$ under the square root.

$$s.d.(\bar{Y}) = \mathbf{S}\sqrt{\frac{1}{n} - \frac{1}{N}} \equiv \sqrt{(1-f)}\mathbf{S}\sqrt{\frac{1}{n}} \quad \text{----(2.10.13)}$$

$$\hat{s.d.}(\bar{Y}) = s\sqrt{\frac{1}{n} - \frac{1}{N}} \equiv \sqrt{(1-f)}s\sqrt{\frac{1}{n}} \quad \text{----(2.10.14)}$$

$$s.d.(P) = \sqrt{\frac{N}{N-1}}\mathbf{P}\mathbf{Q}\left(\frac{1}{n} - \frac{1}{N}\right) \equiv \sqrt{(1-f)}\frac{N}{N-1}\mathbf{P}\mathbf{Q}\frac{1}{n} \quad \text{----(2.10.4)}$$

$$\hat{s.d.}(P) = \sqrt{\frac{n}{n-1}}pq\left(\frac{1}{n} - \frac{1}{N}\right) \equiv \sqrt{(1-f)}pq\frac{1}{n-1} \quad \text{----(2.10.6)}$$

$$= \sqrt{(1-f)}\frac{n}{n-1}pq\frac{1}{n} \quad \text{----(2.10.15)}$$

$$s.d.(\bar{Y})_{\text{EPSWIR}} = \sqrt{\frac{N-1}{N}}\mathbf{S}\sqrt{\frac{1}{n}} \quad \text{----(2.10.16)}$$

- $(1-f)$ in the five second forms makes *less* clear the relative roles of n and N (see Note 2 near the top of the facing page 2.82) and in equations (2.10.13) and (2.10.14) makes it appear (wrongly) that the (data) s.d.s \mathbf{S} [and s] are equivalent to the (probabilistic) s.d. σ – comparing equations (2.10.12) at the right above and (2.10.16) shows that the *actual* equivalence involves a premultiplier (although its value is typically *very* close to 1) for \mathbf{S} – see also equation (2.3.6) on page 2.40.
- Equation (2.10.16) is derived from the *first* term of equation (2.3.9) near the top of page 2.41 in Figure 2.3 when the *second* covariance term is zero due to probabilistic independence under equiprobable selecting *with* replacement (EPSWIR).
- We note that the *reciprocal* of the premultiplier of \mathbf{S} in equation (2.10.16) occurs in equation (2.10.4) to convert the probabilistic (binomial proportion) s.d. $\mathbf{P}\mathbf{Q}$ to the *data* s.d. needed when selecting *is without* replacement.
- When $n=1$, the RHS of equation (2.10.16) is that of equation (2.3.6) on page 2.40 of Figure 2.3.
- When $n=1$, equations (2.10.13) and (2.10.16) become the *same*, reminding us that selecting *without* and *with* replacement no longer differ for a sample of size one. Regrettably, this common special case of equations (2.10.13) and (2.10.16) has been made part of the (unprofitable) $n-1$ vs n saga – see Statistical Highlights #94 and #100 (the lower half of page HL94.9, the middle of page HL100.8).
- * $1-f$ is called the *finite population correction* (abbreviated f.p.c.); such a name conveys the bizarre misconception that we should think of the real world as a ‘correction’ to a model.
- As these Materials maintain throughout, populations in statistics are real-world entities with *finite* numbers of elements and they should *not* be confused with (sometimes useful) infinite *models*.
- * The idea of an infinite ‘population’ is, in part, a confusion of selecting *without* replacement, which ends when all elements have been selected, and selecting *with* replacement which can, in principle, continue indefinitely. However, as the discussion in Note 7 on the facing page 2.82 illustrates, an urn model has a *finite* number of elements even when selecting *is with* replacement (and so can continue indefinitely) under a binomial model.

3. Confidence Intervals for \mathbf{P} , the Respondent Population Proportion

On the basis of approximate normality of the distribution of P as a consequence of the Central Limit Theorem, equation (2.10.6) developed on page 2.81 leads, reasoning as in

Figure 2.3, to an interval involving random variables:

$$I = [P - z_{\alpha}^*\sqrt{\frac{n}{n-1}}\mathbf{P}\mathbf{Q}\left(\frac{1}{n} - \frac{1}{N}\right), P + z_{\alpha}^*\sqrt{\frac{n}{n-1}}\mathbf{P}\mathbf{Q}\left(\frac{1}{n} - \frac{1}{N}\right)],$$

such that $\Pr(I \ni \mathbf{P}) \approx 100(1 - \alpha)\%$, where z_{α}^* is the $100(1 - \alpha/2)\text{th}$

percentile of the $N(0, 1)$ distribution. For calculating an approximate $100(1 - \alpha)\%$ confidence interval for \mathbf{P} we use:

$$p \pm z_{\alpha}^*\sqrt{\frac{n}{n-1}}pq\left(\frac{1}{n} - \frac{1}{N}\right) = [p - z_{\alpha}^*\sqrt{\frac{n}{n-1}}pq\left(\frac{1}{n} - \frac{1}{N}\right), p + z_{\alpha}^*\sqrt{\frac{n}{n-1}}pq\left(\frac{1}{n} - \frac{1}{N}\right)]. \quad \text{----(2.10.17)}$$

NOTES: 9. The theory developed in Figure 2.3 for averages and in this Figure 2.10 for proportions are similar; a *difference* is the (appropriate approximation for the) distributions of \bar{Y} and P , the random variables representing, respectively, the sample average and the sample proportion under equiprobable selecting.

- As discussed in more detail in Note 12 overleaf on page 2.84, use of the normal (rather than the t_{n-1}) distribu-

NOTES: 9. ● tion in the foregoing expressions (2.10.17) for a confidence interval for \mathbf{P} is based on the normal approximation to the binomial distribution (used to model *counted* data).

- The only approximate normality of the distribution of P is one reason why the confidence interval expressions (2.10.17) are only approximate.
10. The theoretical foundation for the confidence intervals expressions (2.3.18) and (2.3.19) on page 2.42 in Figure 2.3 and hence of (2.10.17) overleaf on page 2.83, involves *explicitly* unit inclusion probabilities based on EPS; there is thus *no* basis for the (not uncommon) use of these expressions to calculate confidence intervals for samples selected by *other* (non-probability) methods (accessibility, haphazard, judgement, quota, systematic, volunteer, etc.).
11. The value of a sample attribute determines two characteristics of the approximate confidence interval for \mathbf{P} – the sample proportion defines both its *centre* and its *width*; *both* characteristics (centre *and* width) of a confidence interval may be adversely affected by *compromised* selecting or measuring processes.
- Estimating the population proportion by the sample proportion in estimating the standard deviation of P is another reason why the confidence interval expressions (2.10.17) are approximate.
12. In principle, the normal approximation involved in the confidence interval expressions (2.10.17) overleaf on page 2.83 could be avoided by using the ‘exact’ hypergeometric model (2.10.10) discussed in Note 7 on page 2.82. Although charts and tables have been published for selected values of \mathbf{N} , n and p (see Cochran, page 57, for references), which allow hypergeometric confidence limits for \mathbf{P} to be found, Barnett (page 44) states that their practical utility is limited by factors such as the accuracy with which the charts can be read and often a need for complicated inverse interpolation; however, these difficulties could now likely be mitigated by suitable computer software.
- Provided that n is much smaller than both \mathbf{R} and $\mathbf{N}-\mathbf{R}$, this hypergeometric distribution can be approximated by a *binomial* distribution with parameters n and \mathbf{R}/\mathbf{N} , denoted $Bi(n, \mathbf{R}/\mathbf{N})$. However, computational difficulties again make it inconvenient to find confidence intervals for \mathbf{P} from this binomial approximation. Cochran (pp. 59, 60) gives examples involving calculations for both the hypergeometric and binomial distributions.
 - For practical convenience, it is usual to take the approximation one stage further to the *normal* approximation to the binomial distribution; this approximation is shown at the right in equation (2.10.18).
- $$P \doteq N[\mathbf{P}, \sqrt{\frac{\mathbf{N}}{\mathbf{N}-1} \mathbf{P}\mathbf{Q}(\frac{1}{n} - \frac{1}{\mathbf{N}})}] \quad \text{----(2.10.18)}$$
- The accuracy of the normal approximation will be reasonable provided that:
 - * $n \ll \mathbf{R}, \mathbf{N}-\mathbf{R}$; *i.e.*, the sample size is much smaller than the numbers of elements in the respondent population with and without the characteristic C ;
 - * $n\mathbf{P}$ and $n\mathbf{Q}$ are not too small (say at least 30), because investigating a characteristic that occurs in very few or in nearly all elements of the population generally needs a very large sample size to estimate its proportion with reasonable precision. Table 2.10.2 at the lower right, taken in part from Cochran (page 58), gives working guidelines for deciding when it is reasonable to use the normal approximation; np is the number of selected units in the *smaller* of the categories defined by C and \bar{C} , and n is the required sample size. Cochran states that the rules by which his table is constructed ensure that, at the 95% confidence level, the true frequency with which the confidence limits fail to cover \mathbf{P} is not greater than 5.5%; also, the probability that the upper limit is below \mathbf{P} is between 2.5% and 3.5%, and the probability that the lower limit exceeds \mathbf{P} is between 2.5% and 1.5%. The last column of Table 2.10.2 shows that, under the rules used to construct the table, the (estimated) mean of the binomial distribution, np , lies between about $5\frac{1}{2}$ and $8\frac{1}{2}$ (estimated) standard deviations above the lower limit (0) of the distribution.
- Thus, the normal approximation, in conjunction with equation (2.10.6) near the bottom of page 2.81 leads to the (approximate) confidence interval expressions (2.10.17) overleaf near the bottom of page 2.83, where \mathbf{P} and \mathbf{Q} are estimated using p and q taken as the sample p and q values.
- When n is near the limit where it is reasonable to use the normal approximation, a slight improvement can be affected by incorporating a **continuity correction** of $1/2n$; the expression for calculating the confidence interval for \mathbf{P} is then equation (2.10.19) shown at the right. Without the continuity correction, the normal approximation usually gives a confidence interval that is too narrow (Cochran, page 58).
- $$p \pm z_{\alpha}^* \left[\sqrt{\frac{n}{n-1} pq \left(\frac{1}{n} - \frac{1}{\mathbf{N}} \right)} + 1/2n \right] \quad \text{----(2.10.19)}$$

13. As indicated previously in Note 3 on page 2.82, the expression for *s.d.(P)* in equation (2.10.4) can be used to assess precision. If $n \ll \mathbf{N}$ so that $1/\mathbf{N}$ is negligible compared to $1/n$ and in the ‘worst-case’ [*i.e.*, maximum *s.d.(P)* \equiv minimum precision] situation where $\mathbf{P} = \mathbf{Q} = 1/2$, Table 2.10.3 at the upper right of the facing page 2.85 shows both

Figure 2.10. UNSTRATIFIED POPULATIONS: One-Stage EPSWOR of Individual Elements Estimating a Proportion or a Frequency (continued 2)

NOTES: 13. the resulting standard deviations and 95% confidence interval half-widths (in %) for selected values of n .
(cont.)

- Numbers such as those in the third and sixth columns are often the basis for the statements used in the media to report the precision of polls. For example, for polls taken before elections in Canada, which typically have around 1,500 respondents, the 'margin of error' is commonly given as plus or minus $2\frac{1}{2}$ percentage points. However, in some instances, *larger* (more conservative?) margins of error are given than those in Table 2.10.3; e.g., plus or minus 4 percentage points for (a sample size of) about 1,000 respondents.

Table 2.10.3

n	$\sqrt{\frac{1}{4n}}$	$196\sqrt{\frac{1}{4n}}$	n	$\sqrt{\frac{1}{4n}}$	$196\sqrt{\frac{1}{4n}}$
50	.0707	13.9%	1,000	.0158	3.1%
100	.05	9.8	1,500	.0129	2.5
200	.0354	6.9	2,000	.0112	2.2
250	.0316	6.2	2,500	.01	2.0
300	.0289	5.7%	3,000	.0091	1.8%
400	.025	4.9	3,500	.0085	1.7
500	.0224	4.4	4,000	.0079	1.5
750	.0183	3.6	5,000	.0071	1.4

Example 2.10.1: A small university has 300 final-year undergraduates; in a poll of 100 of these students obtained by equiprobable selecting, 15 said they intended to go on to graduate school. Find an approximate 95% confidence interval for the proportion of all such students who intend to go on to graduate school.

Solution: We have: $N = 300$, $n = 100$, $p = 0.15 (= p)$, $z_{\alpha}^* = 1.95996$ for 95% confidence,

$$\text{so that: } \hat{s.d.}(P) = \sqrt{\frac{n}{n-1}pq\left(\frac{1}{n} - \frac{1}{N}\right)} = \sqrt{\frac{100}{99}(0.15 \times 0.85)\left(\frac{1}{100} - \frac{1}{300}\right)} = 0.029\ 301\ 635.$$

Hence, an approximate 95% confidence interval for P , the respondent population proportion of all such students who intend to go on to graduate school, is:

$$p \pm 1.95996 \times \hat{s.d.}(P) = 0.15 \pm 1.95996 \times 0.029\ 301\ 6 \Rightarrow (0.0926, 0.2074) \text{ or about } (9, 21\%).$$

NOTES: 14. Despite the high sampling fraction of one third, the confidence interval is relatively *wide* – i.e., the small sample size of 100 means that sample error imposes appreciable limitation on the answer involving a (fairly small) respondent population proportion of (we estimate) around 15%.

15. With a sample size of 100, the continuity correction of $\pm 1/2n$ is ± 0.005 ; the confidence intervals without and with the correction are given in Table 2.10.4 at the right below. We see that, in this situation, the difference is unlikely to be practically important – we would typically give a final value for both intervals as about (9, 21%).

Table 2.10.4: Confidence interval		Interval width	Interval width $\div p$
Calculated	Value		
without continuity correction	(0.0926, 0.2074)	0.1148	77%
with continuity correction	(0.0876, 0.2124)	0.1248	83%

Example 2.10.2: A sociological survey is carried out in a small town containing 621 households; 60 of the households obtained by equiprobable selecting show that 11 contain at least one person over the age of 65. Find an approximate 90% confidence interval for the proportion of households in the town that contain at least one person over the age of 65.

Solution: We have: $N = 621$, $n = 60$, $p = \frac{11}{60} = 0.18\dot{3} (= p)$, $z_{\alpha}^* = 1.64485$ for 90% confidence,

$$\text{so that: } \hat{s.d.}(P) = \sqrt{\frac{n}{n-1}pq\left(\frac{1}{n} - \frac{1}{N}\right)} = \sqrt{\frac{60}{59}\left(\frac{11}{60} \times \frac{49}{60}\right)\left(\frac{1}{60} - \frac{1}{621}\right)} = 0.047\ 879\ 847.$$

Hence, an approximate 90% confidence interval for P , the respondent population proportion of households with at least one person over the age of 65, is:

$$p \pm 1.64485 \times \hat{s.d.}(P) = 0.18\dot{3} \pm 1.64485 \times 0.047\ 879\ 8 \Rightarrow (0.1046, 0.2621) \text{ or about } (10, 27\%).$$

NOTES: 16. Despite the appreciable sampling fraction of about ten per cent, the confidence interval is, as in Example 2.10.1, relatively *wide* – i.e., the small sample size of 60 means that sample error imposes appreciable limitation on the answer involving a (fairly small) respondent population proportion of (we estimate) around 18%.

17. With a sample size of 60, the continuity correction of $\pm 1/2n$ is ± 0.0083 ; the confidence intervals without and with the correction are given in Table 2.10.5 at the right. We see that, in this situation, the difference is unlikely to be practically important – we would typically give a final value for both intervals as about (10, 27%).

Table 2.10.5: Confidence interval		Interval width	Interval width $\div p$
Calculated	Value		
without continuity correction	(0.1046, 0.2621)	0.1575	86%
with continuity correction	(0.0963, 0.2704)	0.1741	95%

(continued overleaf)

4. Confidence Intervals for \mathbf{NP} , the Respondent Population Number or Frequency

Under the assumption that the population size, \mathbf{N} , is a *known constant*, the theory of equiprobable selecting for estimating a population number or frequency is a straight-forward extension of that for a proportion. Because the population frequency is \mathbf{NP} , its estimator is \mathbf{NP} ; the standard deviation

of this estimator is then $\mathbf{N} \times s.d.(P)$. Hence,

$$I = [\mathbf{NP} - z_{\alpha}^* \sqrt{\frac{n}{n-1} \mathbf{N}^2 P Q (\frac{1}{n} - \frac{1}{\mathbf{N}})}, \mathbf{NP} + z_{\alpha}^* \sqrt{\frac{n}{n-1} \mathbf{N}^2 P Q (\frac{1}{n} - \frac{1}{\mathbf{N}})}],$$

we obtain an interval involving random variables:

such that $\Pr(I \ni \mathbf{NP}) \approx 100(1 - \alpha)\%$, where z_{α}^* is the $100(1 - \alpha/2)$ th percentile of the $N(0, 1)$ distribution. For calculating an approximate $100(1 - \alpha)\%$ confidence interval for \mathbf{NP} , we use:

$$\mathbf{NP} \pm z_{\alpha}^* \sqrt{\frac{n}{n-1} \mathbf{N}^2 p q (\frac{1}{n} - \frac{1}{\mathbf{N}})} = [\mathbf{NP} - z_{\alpha}^* \sqrt{\frac{n}{n-1} \mathbf{N}^2 p q (\frac{1}{n} - \frac{1}{\mathbf{N}})}, \mathbf{NP} + z_{\alpha}^* \sqrt{\frac{n}{n-1} \mathbf{N}^2 p q (\frac{1}{n} - \frac{1}{\mathbf{N}})}]. \quad \text{----(2.10.19)}$$

The discussion near the end of Note 12 on page 2.84 shows that including a continuity correction of $\mathbf{N}/2n$ in the confidence interval expressions (2.10.19) may be desirable when the sample size is near the limit for reasonable use of the normal approximation to the binomial distribution.

Example 2.10.3: An auditor obtains, by equiprobable selecting, 40 accounts from 1,200 accounts of a certain firm; she finds that, for 16 of the 40 accounts, the underlying documentation is not in compliance with stated procedures. Find an interval estimate of the *number* of the firm's accounts not in compliance.

Solution: We have: $\mathbf{N} = 1,200$, $n = 40$, $p = \frac{16}{40} = 0.4 (= p)$, $z_{\alpha}^* = 1.95996$ for 95% confidence,

$$\text{so that: } \hat{s.d.}(\mathbf{NP}) = \sqrt{\frac{n}{n-1} \mathbf{N}^2 p q (\frac{1}{n} - \frac{1}{\mathbf{N}})} = \sqrt{\frac{40}{39} (1,200^2 \times 0.4 \times 0.6) (\frac{1}{40} - \frac{1}{1,200})} = 92.553 \ 518.82$$

Hence, an approximate 95% confidence interval for \mathbf{NP} , the number of the firm's accounts not in compliance with stated procedures, is:

$$\mathbf{NP} \pm 1.95996 \times \hat{s.d.}(\mathbf{NP}) = 1,200 \times 0.4 \pm 1.95996 \times 92.553 \ 519 = 480 \pm 181.401 \ 195 \Rightarrow (298.60, 661.40) \\ \text{or about } (290, 670) \text{ accounts.}$$

NOTES: 18. When the confidence *level* is not specified in the question, we take the default as 95%. Also, the two significant figures given for the rounded end points of the final interval reflect the fact that there are only two significant digits in the number of accounts in the sample which were not in compliance with stated procedures.

19. The confidence interval is relatively *wide* – i.e., the small sample size of 40 accounts means that sample error imposes appreciable limitation on the answer, despite a (statistically favourable) respondent population proportion of (we estimate) around 40%.

20. With a sample size of 40, the continuity correction of $\mathbf{N}/2n$ is ± 15 ; the confidence intervals without and with the correction are given in Table 2.10.6 at the right below. We see that, with such a small sample size, the difference in the intervals may be large enough to be of practical importance.

Table 2.10.6: Confidence interval		Interval width	Interval width $\div \mathbf{NP}$
Calculated	Value		
without continuity correction	(298, 662)	364	76%
with continuity correction	(283, 677)	394	82%

Example 2.10.4: Inspection of 200 items obtained by equiprobable selecting from a production line yields 3 defectives.

- If the line produces 10,000 items a day, find an approximate 99% confidence interval for the number of defectives produced in a day.
- Comment briefly on the use of such a sample survey to measure the defect rate of this production process.

Solution: (a) We have: $\mathbf{N} = 10,000$, $n = 200$, $p = \frac{3}{200} = 0.015 (= p)$, $z_{\alpha}^* = 2.57583$ for 99% confidence,

$$\text{so that: } \hat{s.d.}(\mathbf{NP}) = \sqrt{\frac{n}{n-1} \mathbf{N}^2 p q (\frac{1}{n} - \frac{1}{\mathbf{N}})} = \sqrt{\frac{200}{199} (10,000^2 \times 0.015 \times 0.985) (\frac{1}{200} - \frac{1}{10,000})} = 85.300 \ 238.29$$

Hence, an approximate 99% confidence interval for \mathbf{NP} , the number of defectives produced in a day, is:

$$\mathbf{NP} \pm 2.57583 \times \hat{s.d.}(\mathbf{NP}) = 10,000 \times 0.015 \pm 2.57583 \times 85.3002 = 150 \pm 219.719 \Rightarrow (-69.719, 369.719) \\ \text{or about } (0, 370) \text{ defectives.}$$

- The confidence interval for \mathbf{NP} , even in its rounded final form of (0, 370) defectives, has a width nearly $2\frac{1}{2}$ times the magnitude of \mathbf{NP} , reflecting the *severe* limitation imposed on the answer by sample error. Example 2.10.4 again reminds us that:

- appreciable (i.e., resource-intensive) sample sizes may be needed to estimate a (small) proportion with adequate precision – that is, to make acceptable in the Question context the limitation imposed on an answer by sample error;

(continued)

Figure 2.10. UNSTRATIFIED POPULATIONS: One-Stage EPSWOR of Individual Elements Estimating a Proportion or a Frequency **(continued 3)**

Solution: (b) ● the *smaller* the magnitude of the population proportion \mathbf{P} , the *larger* the sample size needed to attain a given level of precision specified *relative to* the magnitude of \mathbf{P}

(cont.)

The typical sample size of around 1,500 respondents in national polls is used in part to provide what is deemed to be an acceptable width (*viz.*, about 2.5 percentage points) for a ‘worst case’ 95% confidence interval and in part to give still acceptable (although lower) precision for poll results for subdivisions of the national data such as males and females or regions (*e.g.*, in Canada, Atlantic Canada, Quebec, Ontario, the Prairie provinces, British Columbia) – recall also Table 2.10.3 at the upper right of page 2.85.

NOTES: 21. The unrealizable *negative* lower end point for the confidence interval above has been rounded to zero (the nearest realizable value) for the final statement of the interval. The negative value reminds us to assess:

- the *precision* of the estimate of \mathbf{NP} .
- the accuracy of the hypergeometric-binomial-normal approximation sequence – see Note 7 on page 2.82.

22. With a sample size of 200, the continuity correction of $\pm N/2n$ is ± 25 ; the confidence intervals with-
out and with the correction are given in Table 2.10.7 at the right. We see that the difference in the intervals might be large enough to be of practical importance.

Table 2.10.7: Confidence interval		Interval width	Interval width $\div \mathbf{Np}$
Calculated	Value		
without continuity correction	(0, 370)	370	247%
with continuity correction	(0, 395)	395	263%

23. A summary of information relevant to precision from the four foregoing Examples is given in Table 2.10.8 at the right; the Examples are ordered by sample size. The last column is the estimated *coefficient of variation*, given by $\hat{s.d.}(P)/p$ [equation (2.10.9)], a measure of *relative precision* of the estimator P of \mathbf{P} ; the standard deviation is a measure of *absolute* precision. These coefficient of variation values are of interest in light of the comment by Dr. Glenn in the left-hand column near the top of page 2.50 in Figure 2.3 about a c.v. of no more than 30% being needed for a proportion to be estimated with acceptable sample error.

Table 2.10.8				
Example	n	p	$\hat{s.d.}(P)$	$\hat{c.v.}(P)$
2.10.3	40	0.4	0.0771	19%
2.10.2	60	0.183	0.0479	26%
2.10.1	100	0.15	0.0293	20%
2.10.4	200	0.015	0.00853	57%

NOTES: 24. The theory in this Figure 2.10 for estimating a proportion in two categories under EPS involves *four* standard deviations – two data s.d.s \mathbf{S} and s , two probabilistic s.d.s $s.d.(P)$ and $\hat{s.d.}(P)$ [the last one is the *estimated* standard deviation of P (when \mathbf{P} and \mathbf{Q} are estimated by p and q)]; recognizing and maintaining the distinctions among these s.d.s [and the (model) random variable S and its value s] is part of understanding and applying the theory.

25. The standard deviation of P is sometimes called the *standard error* of P (*e.g.*, Barnett, p. 45); while you should recognize this usage when you encounter it, it has been avoided in these Course Materials because the term is not always used with the same meaning by different authors (see also Cochran, pages 24, 25-27 and 53).

26. A more general reminder about terminology concerns three other terms:

- when you encounter the term *standard deviation* or *standard error*, be sure you understand whether it refers to:
 - the variation of data (and whether the data are from a *sample* or a *census*), **OR**
 - a random variable [and whether it is an individual random variable or a (linear) combination (*e.g.*, an average, sum or difference)], **OR**
 - a parameter of a response model or probability model.
- when you encounter the word *bias*, remember that it is defined only in the context of *repetition* of a *process* – for instance, selecting or measuring.
 - *Estimating* bias *differs* from, for example, (real-world) measuring *inaccuracy* in that it *decreases* with *increasing* sample size and so may not be of much practical concern in actual sample surveys.

Recall that bias (and rms error) are discussed in Appendix 5, Appendix 6 and Appendix 7 on pages 2.46 to 2.49 in Figure 2.3.

5. Appendix 1: Terminology Review [optional reading]

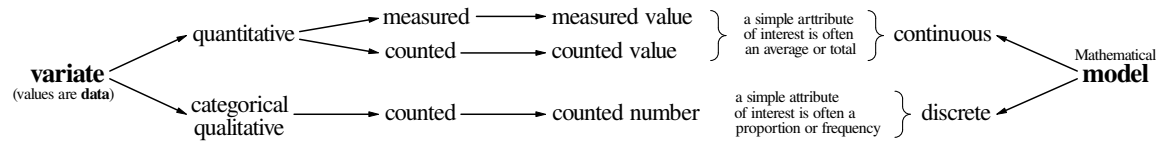
* **Variate:** a characteristic associated with each *element* of a population/process.

- **Response variate (\mathbf{Y}):** a variate defined in the Formulation stage of the FDEAC cycle; an Answer gives some attribute(s) of the response variate over the target population/process.
- **Explanatory variate (\mathbf{Z}):** a variate, defined in the Formulation stage of the FDEAC cycle, that accounts, at least in part, for changes from element to element in the value of a response variate.

(continued overleaf)

– **Focal (explanatory) variate (X):** for a Question with a *causative* aspect, the focal variate is the *explanatory* variate whose relationship to the *response* variate is involved in the Answer(s) to the Question(s).

Distinctions associated with *other* adjectives which qualify 'variate' – **quantitative**, **categorical** (or **qualitative**), **measured** and **counted**, **continuous** and **discrete** – are shown in the following schema.



- An illustration of the distinction between a counted *value* [a real number, modelled by a **continuous** variate] and a counted *number* [an integer, modelled by a **discrete** variate] is:

- to assess the effectiveness of an insecticide, the number of insects could be counted on a defined part of each of a sample of plants from *untreated* and *treated* crops – a measure of effectiveness would be the decrease in the *average* number of insects per plant;
- to assess gender balance in an area of employment, the number of men and women employed in the area could be counted – a (simple) measure of the balance would be the *proportion* of each sex in the area.

The quantitative *measured* and quantitative *counted* distinction blurs progressively as counts become larger in magnitude; it is also affected by the limited resolving power (finite precision) of real measuring processes.

Continuous and discrete variates are *model* concepts because real measuring instruments with finite precision yield only *discrete* values.

- *Quantitative* variate values can become (ordinal) *categorical* – e.g., ages can be classified into age *groups*; we take *qualitative* to mean nominal (*non-ordinal*) *categorical* – e.g., marital status or skin colour.
- A *binary* variate is a categorical variate in *two* categories.

6. Appendix 2: Notation Review [optional reading]

Random variables (discrete or continuous) are a component of the probability *models* used in statistics; *our* concern is more commonly with *discrete* random variables which take on a *finite* number of values (which are often numerical but may be categorical). We denote random variables by upper-case *italic* letters, usually from near the end of the alphabet; we try to keep upper-case italic letters near the start of the alphabet for *events* (bearing in mind the choice of a letter which needs to be *evocative*).

- Formally, a random variable is an entity that assigns a probability to each point of the sample space.
- Informally, a random variable is an entity that takes on different values according to chance.

Both statements mean that we can specify a random variable by a table that gives its values and their probabilities. For example, for selecting equiprobably two balls from an urn containing three red balls and three black balls, the probability functions for the number of red balls in the sample are shown at the right for selecting without and with replacement.

r	0	1	2
$f(r)$	$\frac{1}{5}$	$\frac{3}{5}$	$\frac{1}{5}$

r	0	1	2
$f(r)$	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$

The notation in these tables is a lower case *italic* letter (r) for the random variable *values* and f for the **probability function** that gives the *probabilities*. Other examples of probability functions are the hypergeometric and binomial p.f.s in equations (2.10.10) and (2.10.11) on the lower half of page 2.82 and the p.f. of the random variable V_j in equation (2.10.1) near the middle of page 2.81.

Our use of f (rather than p) for a probability function allows us to keep the letter p for a *proportion*, with four meanings as shown at the right in Table 2.10.9. The distinctions we can maintain with this notation are easily lost when the predominant use of letter p is P for probability, as is the usual practice in introductory statistics texts. In particular, our distinguishing p and P reminds us of the statistical requirement for the Plan for an investigation to make it reasonable to treat sample proportions as values of a (model) random variable P .

Table 2.10.9	
QUANTITY	SYMBOL
Population proportion	P
Sample proportion	p
Random variable	P
A value of P	p
Model parameter	π
Probability	Pr

7. REFERENCES: 1. Barnett, V. *Sample Survey Principles and Methods*. Second edition, Edward Arnold, London, 1991; (First edition: *Elements of Sampling Theory*. The English Universities Press Ltd., London, 1974).

2. Cochran, W.G. *Sampling Techniques*. John Wiley & Sons, Inc., New York, 3rd Edition, 1977.

Appendix 8 of Figure 2.3 continued from page 2.50

It is clear from a few quick calculations that all of the "answers" cited on page 2.50 are based on *two* observations. One was the white female in Region 2, the other the black male in Region 3, both of whom admitted to being regular users of heroin. These were "projected" into 1,440 and 1,320 respectively, and the rest was simple arithmetic. However,

arithmetic cannot erase the very tenuous nature of the answers. The only valid answer is that, in the general population of North Carolina, two people were ready to admit to the regular use of heroin. One was a white female in Region 2; the other was a black male in Region 3. What does this tell us about regular users of heroin in the state?

Very little, except that out of 2,007 people interviewed, two were willing to admit being regular users. Apparently one person in 1,000 interviewed is willing to admit to regular use of heroin. The number who are *actually* regular users and their sociodemographic characteristics remain as unknown as they were before the survey.