University of Waterloo                                                                 STAT 332 – W. H. Cherry

# Assignment 3

**A3 – 1.** Some populations (*e.g.*, the records of the number of days absent over a specified time period for each of the employ-ees of a company or organization) consist of a proportion $P$ of elements with non-zero $Y$-values and a proportion $Q$ ($Q = 1 - P$) of elements with zero $Y$-values. It is sometimes possible to identify accurately and cheaply the $NQ$ ele-ments of the latter type and, if it is desired to estimate $\overline{Y}$ (the population average) from a sample survey, it may then be more efficient to obtain, by equiprobable selecting (EPS), n units (or elements) from the $N_{nz} = NP$ elements with *non*-zero values, rather than from the whole respondent (or study) population of $N$ elements.

(a) If $\overline{Y}_{nz}$ is the average of the $NP$ population elements with non-zero $Y$-values, show that $\overline{Y} = P\overline{Y}_{nz}$.

(b) If $S_{nz}$ is the (data) standard deviation of the $NP$ population elements with non-zero $Y$-values, establish the result at the right, where $S$ is the (whole) population (data) standard deviation.
$$S_{nz} = \sqrt{\frac{N_{nz}}{N_{nz}-1} \left[ \frac{N-1}{N} \cdot \frac{S^2}{P} - \frac{Q}{P^2} \overline{Y}^2 \right]}.$$

(c) If $N_{nz} \gg 1$, write the expression in (b) for $S_{nz}$ in *simpler but approximate* form.

(d) Suppose that $\overline{y}_{nz}$ is the average of n units, obtained by equiprobable selecting, with non-zero $Y$-values. Suggest an appropriate estimator of $\overline{Y}$ based on $\overline{y}_{nz}$, and give its standard deviation in terms of $N$, n, $P$, $Q$, $\overline{Y}$ and $S$. Indi-cate briefly how this estimator is an improvement over the average of n units, obtained by equiprobable selecting, from *all* $N$ respondent (or study) population units (or elements).

**A3 – 2.** *The Globe and Mail* of September 16, 1993, reported the results of a poll which asked the question: *Which leader do you think would make the best Prime Minister of Canada?* Of those who responded, 39% said Kim Campbell and 19% said Jean Chretien. The report contained a common newspaper statement of survey precision: *The results are based on 1,446 personal interviews and are considered accurate within 2.6 percentage points, 19 times in 20.*

(a) Show how the figure of '2.6 percentage points' is obtained and explain briefly what it means.

(b) Would it be correct to say instead '2.6 percent'? Explain briefly.

(c) In assessing the limitations imposed by error of the results of the poll, a student wrote: *The article does not say what economic class the respondents were from or how old they were, but it does not really have to because the poll is accurate and right 19 times in 20.* Comment.

(d) The student also wrote: *The article does not tell us how many people refused the interview but it seems it does not matter – the people who were interviewed and their age and economic status must be positively correlated with regard to the final results because the poll is 'on the nose' 19 times in 20.* Comment.

**A3 – 3.** A dentist is interested in assessing the effectiveness of a new toothpaste, and so she obtains permission to involve a group of $N = 1,000$ school children in an investigation. Pre-investigation records of *all* the children show an average of 2.2 cavities every six months. After three months of the investigation, the dentist obtains, by equiprobable select-

| Child (*i*) | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| Number of cavities in 3 months ($y_i$) | 0 | 4 | 2 | 3 | 2 | 0 | 3 | 4 | 1 | 1 |

ing, 10 children from the group, and obtains the results tabulated above. From an analysis of these data, answer the dentist's question about the toothpaste; show your calculations and reasoning clearly.

**A3 – 4.** One hundred households, obtained by equiprobable selecting from a village with $N = 850$ households, showed 30 households possessing transistor radios; further, a com-plete record of the number of members ($Y$) in these 30 households gave the data shown above.
$$\sum_{j=1}^{n} y_j = 140 \quad , \quad \sum_{j=1}^{n} y_j^2 = 976$$

(a) Estimate both the proportion and number of households in the village which possess transistor radios, and give an approximate 99% confidence interval for each estimate.

(b) Estimate both the average size, and its standard deviation, of households in the village possessing transistor radios. What assumption about one of the estimates from (a) is needed in the calculations in (b)?

**A3 – 5.** This question refers to the article EM9514 *Soveeignty a winner, poll finds, if Quebec gets link to Canada*, which is re-printed in Figure 2.2c on page 2.17 of the Course Materials.

In the middle and right-hand columns on the first side of Figure 2.2c, two non-overlapping 'half-samples' of 501 people are mentioned, with those in the first half-sample being asked about sovereignty and those in the second asked about sovereignty-association. From the information given in the article, find approximate 90% confidence intervals for:

(a) the proportion of Quebeckers in the respondent population who are in favour of sovereignty;

(b) the proportion of Quebeckers in the respondent population who are in favour of sovereignty-association;

(c) the *difference* [(b) − (a)] in the two proportions in (a) and (b).

(d) Explain briefly whether it is a reasonable Answer from this poll that a greater proportion of Quebeckers support sovereignty-association than support sovereignty.

**A3 – 6.** Suppose that there are more than two candidates for an election. Candidate $A$ has the support of a proportion $P_A$ of decided voters, while candidate $B$ has the support of a proportion $P_B$. We want to estimate the *difference* $P_A - P_B$ and to find a confidence interval for this difference.

(a) Define:
$$V_i = \begin{cases} 1 & \text{if decided voter i intends to vote for candidate } A, \\ -1 & \text{if decided voter i intends to vote for candidate } B, \\ 0 & \text{if decided voter i intends to vote otherwise,} \end{cases}$$

for each of the $N$ decided voters in the population. What population attribute of $V_1, V_2, ....., V_N$ is equal to $P_A - P_B$? Justify your answer.

(b) For $V_i$ defined in (a), write the standard deviation $S$ (at the right) in terms of $P_A$ and $P_B$. $\qquad S = \sqrt{\dfrac{1}{N-1}\sum_{i=1}^{N}(V_i - \overline{V})^2}$

(c) Suggest an estimator of $P_A - P_B$. How does it relate to $V_1, V_2, ....., V_n$, random variables representing the *sample* analogues of the population $V_i$s?

(d) Based on your answers to (a) and (b), give an expression for the standard deviation of the estimator in (c).

(e) Give an expression for the *estimated* standard deviation corresponding to (d).

(f) Suppose that the Liberals have the support of 36% of 2,500 decided voters, obtained by equiprobable selecting, while the Conservatives have 31% support in this sample. Find an approximate 95% confidence interval for the Liberals' lead over the Conservatives.

(g) Comment briefly on the interpretation of the confidence interval found in (f).

1995-04-20