

Statistics 257-001 Applied Sampling Techniques
Final Examination (December 17, 2004)

1. (6 points) Carefully define what is meant by a *simple random sample of size n from a population of size N* .
2. (12 points) List three principal reasons for choosing to use stratified random sampling rather than simple random sampling.
3. (8 points) Carefully discuss two situations in which cluster sampling is an effective design for obtaining a specified amount of information at minimum cost.
4. (12 points) For the following survey situation, *and in the context of the situation described*, carefully state the target population, the frame, and the sampling units. Also discuss any possible sources of selection bias or inaccuracy of responses, if appropriate.

The Saskatchewan Provincial Travel Committee commissioned a study to identify inter-provincial (within Saskatchewan) travel patterns of Regina and Saskatoon residents, and to evaluate different sources of vacation planning information. They conducted 400 interviews with Regina residents and 400 interviews with Saskatoon residents. Telephone numbers with Regina and Saskatoon exchanges were generated randomly so that listed and unlisted telephone numbers could be reached. “Respondents were limited to heads of household and quotas were established in order to have an equal representation of male and female respondents. Additionally, income and age brackets were monitored in order to maintain the same proportions as the general population bases of Regina and Saskatoon.”

5. (18 points) A nutritionist at the University of Saskatchewan has decided he would like to know how much pizza students can eat. Using an official list of registered full-time undergraduate students from the Registrar, he selects a simple random sample of size $n = 17$. He provides pizza at lunchtime for these 17 students, and records the total number of slices each student ate. He decides that he would like to estimate \bar{Y} , the average number of pizza slices that a University of Saskatchewan student can eat at lunchtime, except that he cannot decide which estimator he should use. Let y_i denote the number of pizza slices that the i^{th} student ate, and consider the following two estimators of \bar{Y} :

$$\bar{y} = \frac{y_1 + y_2 + y_3 + \cdots + y_{15} + y_{16} + y_{17}}{17}$$

and

$$\hat{y} = \frac{y_1 + 3y_2 + 3y_3 + \cdots + 3y_{15} + 3y_{16} + y_{17}}{47}.$$

(a) Show that BOTH \bar{y} and \hat{y} are each unbiased estimators of \bar{Y} . This shows that unbiased estimators are not unique!

(b) It is also known from previous pizza experiments that $\text{Var}(y_i) = 2$, $i = 1, \dots, 17$. Compute both $\text{Var}(\hat{y})$ and $\text{Var}(\bar{y})$ exactly stating any assumptions that you have made, and therefore show that in this case $\text{Var}(\hat{y}) \geq \text{Var}(\bar{y})$.

6. (26 points) After hearing of the nutritionist's experiment in Problem 5, a sociologist at the University of Saskatchewan decided to find out if a drug he invented, called *Pizza-X*, could help university students to eat more pizza. He decided to test his drug on his Sociology 101 class. Fortunately, all 100 students in the class agreed in advance to participate if selected for the experiment. Using a table of random digits, he randomly selected 10 students, and allocated 5 to the control group and 5 to the injection group. All participants were instructed not to eat after 10 a.m. and were then fed as many standard-sized cheese pizza slices during lunchtime as could be eaten. The results obtained are listed below:

control group	injection group
11	15
12	11
11	14
7	9
9	11

Let \bar{y}_1 denote the mean number of pizza slices eaten by the control group, and let \bar{y}_2 denote the mean number of pizza slices eaten by the injection group.

(a) Construct an approximate 95% confidence interval for \bar{Y}_1 , the true mean number of pizza slices eaten by those who have not taken *Pizza-X*.

(b) Construct an approximate 95% confidence interval for \bar{Y}_2 , the true mean number of pizza slices eaten by those who have taken *Pizza-X*.

(c) Based on your answers to (a) and (b), is there sufficient evidence to conclude that the sociologist's drug *Pizza-X* helps his Sociology 101 students eat more pizza? Why or why not?

7. (14 points) A sociologist at the University of Regina is interested in the *total* number of families that rent their homes in the City of Moose Jaw. She decides to conduct a 1-in-50 systematic survey drawing upon city tax records. It is known that there are 15 200 households in Moose Jaw. The sociologist found that

$$\sum_{i=1}^{304} y_i = 76$$

where $y_i = 1$ if the family in the i^{th} household sampled rents and $y_i = 0$ if the i^{th} family does not. Using the data given, construct an approximate 95% confidence interval for Y_T , the total number of families who rent.

(Hint: If \bar{y} denotes the estimated *proportion* that rent, then $y_T = N\bar{y}$.)

8. (20 points) The City of Regina chief engineer is considering a zoning change for the Wascana Park subdivision to allow a new shopping complex to be built. In order to assess the opinion of residents, a cluster sample is used. The subdivision map is marked into $M = 170$ blocks, and a random sample of $m = 15$ blocks is surveyed. For each of these blocks (labelled i , for $i = 1, \dots, 15$), the number of adult residents (n_i) and the number of adult residents opposing the zoning change (y_i) are recorded. The data are summarized below:

$$\sum_{i=1}^{15} n_i = 546, \quad \sum_{i=1}^{15} n_i^2 = 9981, \quad \sum_{i=1}^{15} n_i y_{iT} = 4035, \quad \sum_{i=1}^{15} y_{iT} = 182, \quad \sum_{i=1}^{15} y_{iT}^2 = 1819,$$

$$\sum_{i=1}^{15} \bar{y}_i^2 = 1.54, \quad \sum_{i=1}^{15} n_i^2 \bar{y}_i^2 = 2103, \quad \sum_{i=1}^{15} n_i^2 \bar{y}_i = 4571.$$

Use an appropriate estimator to construct an approximate 95% confidence interval for \bar{Y} , the proportion of Wascana Park adult residents opposed to the zoning change.

9. (24 points) A psychologist at the University of Regina is interested in educational reform and wants to obtain information on the amount of time spent on non-academic activities such as sports and video games by elementary school children. She uses stratified sampling as indicated in the table below, collecting information on Y , the percentage of time spent on non-academic activities.

STRATA	strata size	sample size	sample mean	sample variance
Separate Schools	24	6	13	9
Public Schools	54	12	26	16

(a) Find an approximate 95% confidence interval for the population mean \bar{Y} .

Suppose that the psychologist had wanted the variance for the estimation of \bar{Y} to be $V = \text{Var}(\bar{y}_{ST}) = \frac{9}{16}$. What sample size should she have used and how should it have been allocated for each of the allocation methods given in (b) and (c) below?

(b) *proportional allocation:*

(c) *Neyman allocation:*

10. (24 points) Each person in a population of adults is interviewed, an identification number is assigned to each one, and the sex and age of each individual is recorded. The population size is 1000 of which 550 are female and 450 are male. The population mean age is 45. From the list of identification numbers a simple random sample of size 50 is obtained and the following measurements are recorded: x_i is the age (in years) and y_i is the systolic blood pressure (in mg/100ml) of the i^{th} subject.

The resulting data are

$$\sum_{i=1}^{50} x_i = 2248, \quad \sum_{i=1}^{50} x_i^2 = 104384, \quad \sum_{i=1}^{50} y_i = 6744, \quad \sum_{i=1}^{50} y_i^2 = 928436, \quad \sum_{i=1}^{50} x_i y_i = 305125.$$

Find an approximate 95% confidence interval for the population mean of y using the estimation methods given in (a) and (b) below.

(a) *ratio estimation:*

(b) *regression estimation:*

(c) Compute the relative efficiency of these two estimators. Is there sufficient evidence to favour one over the other? Why or why not?

11. (18 points) A senior administrator at the University of Regina wishes to estimate the proportion of its students that have used cocaine, a sensitive subject. Students were classified into one of two strata—undergraduate and graduate—and were randomly sampled within the stratum. Since there was some concern that students might be unwilling to disclose their use of cocaine to a university official, the following *random response method* was used. The university official constructs a deck of 30 cards. On 26 of them are marked **N** for *never used cocaine* and 4 of them are marked **C** for *have used cocaine at least once*. Each sampled student was asked to draw a card from the deck and to respond *yes* if the letter agrees with the group that student belongs to. The results are as follows:

STRATA	total number of students	number sampled	number answering yes
undergraduate	8972	900	723
graduate	1548	150	117

(a) Construct a 95% confidence interval for the proportion of undergraduates who have used cocaine at least once.

(b) Construct a 95% confidence interval for the proportion of graduates who have used cocaine at least once.

(c) Based on your answers to (a) and (b), is there a statistically significant difference in the use of cocaine between undergraduates and graduates? Is this surprising? Why or why not?

12. (18 points) Suppose that the City of Regina wants to estimate how many homeowners band their elm trees to prevent *Dutch Elm Disease*. A sample is to be selected using one of the following: simple random sampling, systematic sampling, cluster sampling. The method of data collection could be one of: personal interview, telephone survey, mailed questionnaire, direct observation. Choose both a sampling scheme and a method of data collection, and discuss how you might actually conduct the survey using your two choices.

(In answering this question, you might want to draw on many of the ideas discussed throughout the course. If there are any benefits or limitations to your plan, you should mention them as well.)