| | | |
|---|---|---|
| Slide 1 |  Unit 1 Lesson 6<br><br>**Sampling Example**<br>Dr. Jane Monaco<br>Clinical Assistant Professor<br>Department of Biostatistics<br>Gillings School of Global Public Health<br>The University of North Carolina at Chapel Hill | |
| Slide 2 | **Objectives**<br><br>▪ Review Simple Random Sampling, Stratified Sampling and Multi-stage Sampling<br>▪ Understand the role of sampling in motivating example: Quality of Life and Obesity in School Children in Australia | ▪Review Simple Random Sampling, Stratified Sampling and Multi-stage Sampling<br>▪Understand the role of sampling in motivating example: Quality of Life and Obesity in School Children in Australia – one of our global health examples. |
| Slide 3 | **SIMPLE RANDOM SAMPLE**<br><br>A simple random sample (SRS) is a probability sample of size $n$ in which every set of $n$ units has exactly the same chance of being selected<br><br>*Simplest and most important method of sampling* | Recall that a probability sample is one in which every element has a known non-zero probability of selection and that probability is determined by chance. Most important example of a probability sample is a simple random sample or SRS. Many other methods use a SRS as a starting point.<br><br>An SRS is a probability sample of size $n$ in which every set of $n$ units has exactly the same chance of being selected. [This is an important definition to know and understand.]<br><br>An SRS is the simplest and most important method of sampling, because it is a building block we need other methods and when we talk about randomization.<br><br>An SRS is a probability sample, but not the only kind of probability sample. |

| | | |
|---|---|---|
| Slide 4 | **SELECTING AN SRS:**<br>Method 1- GROUPING OF DIGITS<br>• Label units in the population from 1 to $N$<br>• Select a row in random number table<br>• Group the digits by the number of digits in $N$<br>    Example:    $N$=9, group by 1 digit<br>                      $N$=2500, group by 4 digits<br>• Select units with labels of grouped digits, ignore numbers greater than $N$ and ignore duplicates<br><br>UNC GILLINGS SCHOOL OF GLOBAL PUBLIC HEALTH   4 | The first method for selecting an SRS that we'll discuss is to group digits in the random number table.<br>Here's the method in general, and then we'll look at in an example….<br>First, label units from 1 to $N$ where "$N$" is the population size.<br>Then, select a row in random number table available in most introductory statistics text books (or random number book). (Can you imagine a more boring book than a random number book? Yet …there is such a thing…) You can select any row… we'll talk more about that in a minute.<br>Then, group the digits in the table by the number of digits in $N$<br>  [$N$=9, group by 1 digit, $N$=2500 group by 4 digits, $N$= 72 group by 2 digits….]<br>Finally, select units with labels of grouped digits, ignoring numbers greater than $N$ and ignoring duplicates.<br>It sounds complicated when you write down the steps, but we've seen it is pretty straightforward. |
| Slide 5 | **SELECTING AN SRS:**<br>Method 2- Multiplication by a URN<br>• Label units from 1 to $N$.<br>• Select a row in random number table.<br>• Obtain a URN (uniform random number), call it $u$, between 0 and 1.<br>• Multiply $u$ by $N$ (population size).<br>• Select the <u>next larger integer</u><br>• Continue (discarding duplicates) until all $n$ units have been selected.<br><br>UNC GILLINGS SCHOOL OF GLOBAL PUBLIC HEALTH   5 | The next method uses multiplication by the uniform random numbers.<br><br>Like the previous method, you start by numbering the members of the population from 1 to $N$. Next you select a row in the random number table.<br>You then insert a decimal before the first digit, to view the random number as a number between 0 and 1.<br>This value is called a uniform random number (URN).<br>You will multiply the random number (between 0 and 1) by $N$, the population size. Take the next largest integer (in other words, round up). Select this unit of the population.<br><br>Repeat these steps, ignoring duplicates until you have selected $n$ units for the sample.<br>Again, when you just read the steps, it looks confusing and hard, but it is not hard in practice. |
| Slide 6 | **DIFFERENT METHODS, DIFFERENT SRSs**<br><br>   – No <u>one</u> correct SRS<br>   – Depends on the method, the numbering of original units, the row you select from the table, ….<br><br>UNC GILLINGS SCHOOL OF GLOBAL PUBLIC HEALTH   6 | Do we get same simple random samples using these methods? In general, no -<br>There is no one correct SRS. The SRS will depend on, among other things, the way you originally ordered/labeled the population, the method you use, the row you happen to select in random table, ….<br>So there are lots of possible valid SRSs from a given population. |

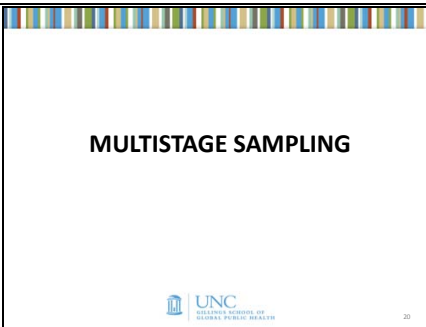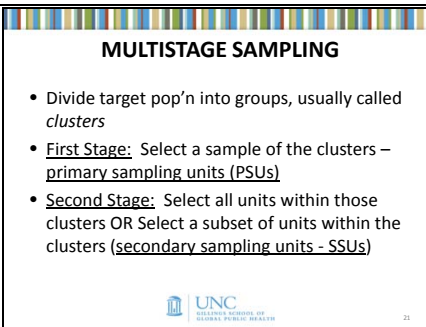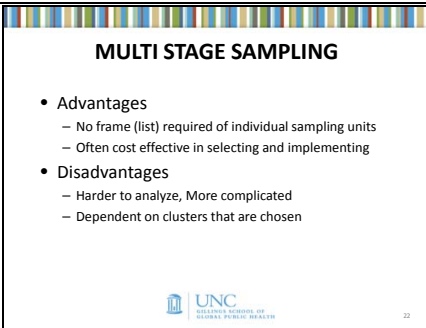| | | |
|---|---|---|
| Slide 7 | **SELECTING AN SRS:**<br>Method 3: Computer Generated<br>• Large SRS or large population → random number table is inconvenient<br>• Uniform random numbers (URNs) can be generated by almost all statistical software packages and many calculators<br>    – Use computer to calculate the uniform random numbers<br>    – Continue as in method 2<br><br>UNC<br>GILLINGS SCHOOL OF<br>GLOBAL PUBLIC HEALTH    7 | The most common method of selecting a SRS in practice is computer-generated.  So instead of using a random number table, you use a computer to generate the random numbers.  This computer method is almost the same as Method 2… In this method, you let the computer or calculator generate the uniform random numbers between 0 and 1.   Any statistical software and lots of calculators have this feature.  Then, you multiply the uniform random number by $N$ ( population size) and take next larger integer, as in method 2.<br>With a little programming, you can set up a loop to select even big sample very quickly. |
| Slide 8 | **MORE COMMENTS ABOUT AN SRS**<br><br>• An SRS requires you to be able to list (enumerate) all the units in the population→ Big disadvantage.<br>• Record the way your random numbers were selected (row or table, seed….)<br>• A SRS may look like it has 'patterns'…. that's OK<br><br>UNC<br>GILLINGS SCHOOL OF<br>GLOBAL PUBLIC HEALTH    8 | A couple of general notes about SRSs:<br>1:  One big disadvantage of the selecting a SRS in general is that you must be able to list all the units in the population.  This is a pretty big disadvantage which will be addressed in future lessons with multi-stage sampling.<br>2: Remember to make sure that you record how you selected your sample- in other words how you found the random number(s).  For example, if you are using the tables, record the row number.  If using the computer, record the "Seed". -  So you can reproduce the sample, if necessary.<br>3:  I noted previously that you can't tell by looking at a sample whether it is selected at random.  This is true.  However, I sometimes do an exercise in a residential class that shows that selection at random may seem like it has a pattern.  I have half the class get together and use one of the methods with a random number table to select a sample.  Then, I ask the other half of the class to  gather and select a sample by selecting sample which "looks" random to them.  They write down IDs so that the IDs look random.  I leave the room while they complete the task.<br>When I return, I try to guess which sample was chosen truly at random (using a table) and which was created to look random.  Usually I am able to do it – the truly random sample may look like it has anomalies… IDs selected in a row (like 5,6,7), or many more even numbers in a row,  or have big gaps,  for example.<br>I recently read a Newsweek article that lots of people think their iPods have favorites – when the iPod is  supposed to shuffle the tunes at random, it may seem repeat one artist or music type while ignoring a particular tune for what seems like forever.   I digress… but the point is… use truly random numbers from a table or computer.  Even if your sample seems to have a pattern, that is OK.  True random sampling, like random shuffling on an iPod, may 'feel' like it plays favorites. |

| | | |
|---|---|---|
| Slide 9 | **STRATIFIED RANDOM SAMPLING**<br><br>UNC GILLINGS SCHOOL OF GLOBAL PUBLIC HEALTH | Sometimes, our study calls for a more complicated sampling scheme than SRS.<br>One such scheme is Stratified Random Sampling. |
| Slide 10 | **STRATIFIED RANDOM SAMPLING**<br><br>• Divide population into strata or groups [groups are similar in some characteristics]<br>  Typical strata: ethnicity, age, geographic region, gender, …<br>• Select a sample within each stratum separately<br>• Combine samples in each stratum for final sample<br><br>UNC GILLINGS SCHOOL OF GLOBAL PUBLIC HEALTH   10 | In a stratified random sample, we divide the population into strata (or groups). These groups are similar in some characteristic… some common examples are to stratify by ethnicity, age, geographic region, etc.<br>Within each stratum, we select a sample (usually a SRS). Then, we combine these samples from each stratum to form our final sample.<br><br>Recall that with a stratified sample, we are not selecting a subset of the strata – we are selecting observations WITHIN each the strata. In other words, if we stratify on ethnicity (5 groups) we don't, say, select 2 out 5 ethnicity, rather we select observations within each of the 5 ethnicities. But we can have some flexibility about how we select within each different ethnicity/strata. |
| Slide 11 | **PROPERTIES OF A STRATIFIED RANDOM SAMPLE**<br>• A stratified random sample is never an SRS.<br>• Common to make stratum sample sizes equal rather than sampling fractions<br><br>UNC GILLINGS SCHOOL OF GLOBAL PUBLIC HEALTH   11 | The properties of a stratified random sample are :<br>A stratified random sample is never SRS. Why?<br>Recall the definition of an SRS… the probability of selection for any sample is the same. With a Stratified Random Sample, there are some samples that are not possible (probability of selection is zero).<br><br>It's common to make stratum sample sizes equal rather than sampling fractions equal. |

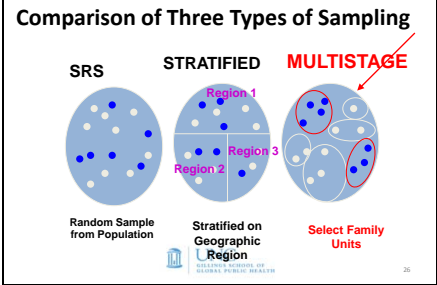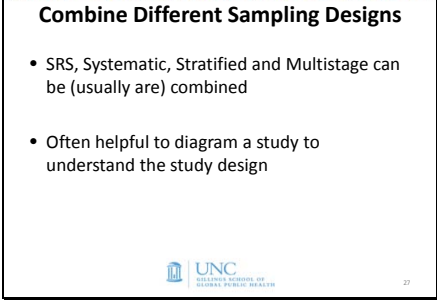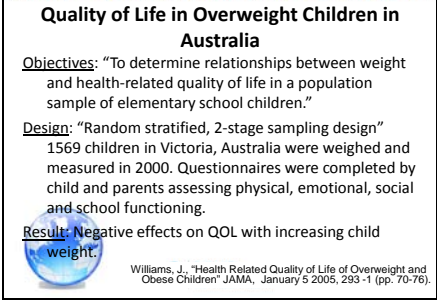| | | |
|---|---|---|
| Slide 12 | **EXAMPLE: STRATIFIED RANDOM SAMPLE**<br><br>• Select 100 patients out of 1000 patients → select equal numbers of diabetics and nondiabetics<br>• Form strata 200 diabetics and 800 nondiabetics<br>• Select 50 diabetics and 50 nondiabetics<br>  – SRS of 50 out of 200 diabetics<br>  – SRS of 50 out of 800 nondiabetics<br>• Full sample is $n$= 100 out of $N$=1000<br><br>UNC GILLINGS SCHOOL OF GLOBAL PUBLIC HEALTH   12 | We saw in an example that we may want to stratify on diabetic status (Diabetic vs. Non-Diabetic)<br>Suppose there are 1000 patients and we wish to select 100 patients for our sample.<br>Suppose the population has 200 diabetics and 800 nondiabetics – these are our two strata.<br><br>Suppose we need to select equal numbers of diabetics and nondiabetics. So select 50 diabetics and 50 nondiabetics<br>Then, select a SRS of 50 out of 200 diabetics and select a SRS of 50 out of 800 nondiabetics. In other words we force that "50/50" split. You can use any method for selecting an SRS.<br>Thus the full sample is n= 100. |
| Slide 13 | **NOTATION: STRATIFIED RANDOM SAMPLE**<br><br>• Divide population into $j$ strata.<br>  – Example has $j$=2 strata (nondiabetics and diabetics)<br>• Number selected in each stratum = $n_j$<br>  – $n_1$=50 diabetics and $n_2$ = 50 nondiabetics.<br>  – Sum of $n_j$'s is $n$<br>• Number in population in stratum $j$ =$N_j$<br>  – $N_1$=200 and $N_2$=800<br>  – Sum of $N_j$'s is $N$<br>• Stratum specific sampling fraction = $n_j/N_j$<br>  – diabetic sampling fraction=$n_1/N_1$=50/200 = 0.25<br>  – non-diabetic sampling fraction $n_2/N_2$ =50/800=0.0625<br><br>UNC GILLINGS SCHOOL OF GLOBAL PUBLIC HEALTH   13 | Recall our notation for stratified samples.<br>Divide population into $j$ strata. (We had $j$=2 strata. Diabetic and non-diabetic)<br>Number selected in each strata is $n_j$. (We selected $n_1$=50 diabetics and $n_2$ = 50 nondiabetics. Convince yourself that the sum of the little nj's are little n)<br>Number in population in stratum $j$ is $N_j$. ($N_1$= 200, $N_2$= 800)<br>The stratum specific sampling fraction is $n_j/N_j$. We just consider the sampling fraction for each group, each stratum. Our example has diabetic sampling fraction=50/200 and nondiabetic sampling fraction 50/800. |
| Slide 14 | **WHY STRATIFY?**<br><br>• Analyze strata separately<br>• If the strata are homogeneous, we can add efficiency<br>• Convenience – different sampling plan for different strata<br><br>UNC GILLINGS SCHOOL OF GLOBAL PUBLIC HEALTH   14 | Why do we stratify?<br><br>One reason, is so that we can analyze strata separately- we can assure that we have sufficient sample size within the strata. We can force by design to have enough observations in a certain strata (diabetic status, ethnicity, etc). (In the previous example we assure that we have enough diabetics to analyze separately.)<br><br>If strata are homogeneous, we can add efficiency – this means that we can get smaller estimates of variance with the stratified sample as compared to using a SRS (under some assumptions). The details go beyond what we cover in Bios 600, but one thing to remember is that stratification can potentially be beneficial not just logistically but can be beneficial by decreasing the necessary sample size or decreasing variability of estimates if the strata are "homogeneous" (meaning similar with respect the outcome of interest).<br><br>We also know that with stratified sampling, we could use not |

| | | only different sampling fractions but also different sampling methods or logistical details.  Perhaps we stratify on age group.  Then not only can we have a different sampling fraction for, say, the older age group but we could have a totally different sampling method….we could use systematic sampling for the older age group and a SRS for the other age group…<br>In general we want a small number of strata that are homogeneous. |
|---|---|---|
| Slide 15 | **OVERSAMPLING AND UNDERSAMPLING**<br><br>• First scenario → selected 50 diabetics and 50 nondiabetics<br>  – <u>Over-sampled</u> diabetics $n_1/N_1 > n/N$<br>  – <u>Under-sampled</u> nondiabetics, $n_2/N_2 < n/N$<br><br>UNC<br>GILLINGS SCHOOL OF GLOBAL PUBLIC HEALTH<br>15 | We have seen it is possible to oversample or undersample the different strata.<br><br>Over-sampling is when the sampling fraction for the stratum is bigger than the overall sampling fraction.<br>Under-sampling is when the sampling fraction for the stratum is smaller than the overall sampling fraction.<br><br>We selected 50 diabetics and 50 non-diabetics… the overall sampling fraction was 0.1.  The sampling fraction for diabetics was $50/200 = 0.25$.  The sampling fraction for the nondiabetics was  $n/N = 50/800= 0.0625$.<br>OK… So the diabetics are over-sampled and the non-diabetics are under-sampled.<br><br>NOTE:  Under-coverage (in a previous lesson) and under-sampling are different concepts. |
| Slide 16 | **EXAMPLE: STRATIFIED RANDOM SAMPLE**<br><br>• Another scenario<br>  $n_1$= 35 diabetics, $n_2$= 65 nondiabetics<br>• Diabetics are oversampled<br>  $n_1/N_1 = 35/200= 0.18 > n/N = 100/1000$<br>• Non-diabetics are undersampled<br>  $n_2/N_2 = 65/800=0.08 < n/N =100/1000$<br><br>UNC<br>GILLINGS SCHOOL OF GLOBAL PUBLIC HEALTH<br>16 | Here is another example, suppose that $n_1$= 35, $n_2$= 65.  This may be because you want to analyze the diabetics separately– the necessary sample size within each stratum would be determined by a statistician to ensure sufficient sample size to analyze diabetics separately.<br>Check that the sampling fraction for the diabetics is bigger than the sampling fraction for the non-diabetics.  So the diabetics are oversampled. |

| Slide 17 |  STRATIFIED RANDOM SAMPLE EXAMPLE (cont)<br><br>• Study investigates the proportion of subjects who achieved a target weight<br>• Let $m_i$ be the number of subjects in the sample within the $i$th stratum to reach a target weight<br>• Suppose $m_1$=10 diabetics out of 35 reach target weight and $m_2$=30 out of 65 non-diabetics reach target weight.<br>• _BIASED_ ESTIMATE IN POPULATION: 40/100 reach target weight. Expect 400 out of 1000 to reach target weight in population? NO. | Suppose that a study investigates the proportion of subjects who achieved a target weight.<br>When we talk about sampling a good bit of time is spent in how to select the sample, but remember eventually we have some dependent variable/ some outcome variable that we want to measure…. And then we want to know not just about the outcome in the <u>sample</u>, but the outcome in the <u>population</u>.<br><br>You'll have some diabetics who reach their target weight, some non-diabetics who reach their target weight and the proportions within the subgroups may be different.  We are still considering the sample to be 35 diabetics and 65 non diabetics.<br>Reviewing the notation…. let $m_i$ be the number of subjects in the sample within the $i$th stratum to reach a target weight.  Suppose $m_1$=10 diabetics out of 35 reach target weight and $m_2$=30 out of 65 non-diabetics reach target weight.<br><br>**BIASED** <u>ESTIMATE OF TARGET WEIGHT FOR POPULATION</u>: Some might say (incorrectly) that therefore we'd expect in general that 40% of people in the population to reach their target weight.  This is because 40 out of 100 reached target weight in the study- in the sample.   But the diabetics made up a larger proportion in the sample than in the population.  So in this <u>biased</u> estimate the diabetics contribute too much to the estimate.  So you need to take into consideration the stratified sampling design. |
| Slide 18 |  UNBIASED ESTIMATE OF PROPORTION<br><br>● Proportion of diabetics reaching target weight:<br>10/35 = 0.286<br>● Proportion of nondiabetics reaching target weight:        30/65 = 0.462<br><br>● Weight these proportions reaching target weight by proportion of diabetics (200/1000) and proportion of nondiabetics (800/1000) in pop'n<br>● <u>Unbiased estimate</u> of proportion reaching target weight in population<br>0.286(0.2) + 0.462(0.8) = 0.0572 + 0.3696<br>= 0.4268      (or about <u>43%</u>) | We can compute the unbiased estimate of proportion  in population reaching target weight.<br> First compute proportion of diabetics reaching target weight in the sample, 10/35 = 0.286  or we could say that 29% of diabetics achieved the target weight in the sample.<br>Proportion of nondiabetics reaching target weight in the sample, 30/65 = 0.462 (or 46%).<br>Now we want to weight these numbers above by proportion of diabetics (200/1000) and proportion of nondiabetics (800/1000) in the population.<br>Unbiased estimate = 0.286(0.2) + 0.462(0.8) = 0.0572 + 0.3696 = 0.4268<br>In other words, in the population of 1000 we would expect that about 43% would achieve their target weight if placed in the same program. If we conducted the study on all 1000 people, we'd expect about 427 of them to achieve their target weight. |

| Slide 19 | **DISADVANTAGES OF STRATIFICATION**<br><br>• Makes analysis more difficult (costly)<br>• May offer no advantages (no appropriate stratification variables)<br><br>UNC GILLINGS SCHOOL OF GLOBAL PUBLIC HEALTH    1919 | Makes analysis more difficult (costly)<br>Stratification may offer no advantages – there may be no variable on which is makes sense to stratify – no variable where the observations are "homogeneous"– or may be expensive measure the stratification variable. |
| --- | --- | --- |
| Slide 20 | **MULTISTAGE SAMPLING**<br><br>UNC GILLINGS SCHOOL OF GLOBAL PUBLIC HEALTH    20 | Another method of sampling that is often used is multistage sampling. Another term for this sampling design is **cluster sampling.** |
| Slide 21 | **MULTISTAGE SAMPLING**<br><br>• Divide target pop'n into groups, usually called *clusters*<br>• First Stage: Select a sample of the clusters – primary sampling units (PSUs)<br>• Second Stage: Select all units within those clusters OR Select a subset of units within the clusters (secondary sampling units - SSUs)<br><br>UNC GILLINGS SCHOOL OF GLOBAL PUBLIC HEALTH    21 | Moving to a different sampling strategy, we will investigate multistage sampling.<br>In multistage sampling, the target population is separated into groups, often called clusters.<br>In the first stage, we randomly select a sample of the clusters (primary sampling units, or PSUs)<br>In the second stage, either we select ALL the units within the clusters which have already been selected OR you may select a subset of units within the clusters. These selections are called secondary sampling units or SSUs. |
| Slide 22 | **MULTI STAGE SAMPLING**<br><br>• Advantages<br>  – No frame (list) required of individual sampling units<br>  – Often cost effective in selecting and implementing<br>• Disadvantages<br>  – Harder to analyze, More complicated<br>  – Dependent on clusters that are chosen<br><br>UNC GILLINGS SCHOOL OF GLOBAL PUBLIC HEALTH    22 | On the positive side, with multistage sampling, no frame (list) required of individual sampling units<br>Also, multistage sampling is often cost effective in selecting and implementing.<br><br>However, multistage design are harder to analyze and just more complicated (expensive).<br>Of course, this method is dependent on which clusters are chosen – so you choose your clustered carefully and in an unbiased way (you do this always, anyway!), but special care needs to be taken in cluster sampling. |

| Slide 23 | **MULTISTAGE SAMPLING: Example**<br><br>• Estimate Proportion of Homeless People who are HIV positive in California<br>  – Select SRS of Counties<br>  – Within each selected County select a homeless shelter<br>  – Within each selected shelter select 5 people to be tested<br><br>UNC KILLINGS SCHOOL OF GLOBAL PUBLIC HEALTH | Estimate Proportion of Homeless People who are HIV positive in California<br>　　Select SRS of Counties<br>　　Within each selected County select a homeless shelter<br>　　Within each selected shelter select 5 people to be tested |
|---|---|---|
| Slide 24 | **Comparison of Three Types of Sampling**<br><br>SRS　STRATIFIED　MULTISTAGE<br><br>Region 1, Region 2, Region 3<br><br>Random Sample from Population　Stratified on Geographic Region　Select Family Units | We've reviewed several types of sampling… let's compare three of the types that we've discussed.<br>Suppose each of these ovals represents the population– copies of the SAME population…the dark blue dots are units within the population that are sampled and the white dots are the units in the population which are not sampled.<br>Each oval represents the **same** population – under the three different sampling strategies, different units are selected.<br>We can think of each dot as a person… and each person is part of a geographic region and a family which may play a part in the sampling scheme.<br>We'll start with SRS.  In a SRS, over to left, we can see our dots which are sampling units.  We just randomly select several sampling units to be our sample- these are the dark blue dots. |
| Slide 25 | **Comparison of Three Types of Sampling**<br><br>SRS　**STRATIFIED**　MULTISTAGE<br><br>Region 1, Region 2, Region 3<br><br>Random Sample from Population　Stratified on Geographic Region　Select Family Units | 　　In a stratified sample, we think of the population as divided into groups called strata.  Here let's say that the data are stratified by geographical location.  There are three strata (say, geographical region 1, 2, and 3).  Then within each strata we select using some method (SRS? Systematic?) individuals within each region.  How many we select within each region  will be determined by the goals of the study.  The strata may be over sampled, under sampled or proportionately sampled. |

| | | |
|---|---|---|
| Slide 26 |  **Comparison of Three Types of Sampling**<br><br>SRS  STRATIFIED  MULTISTAGE<br><br>Region 1<br>Region 3<br>Region 2<br><br>Random Sample from Population  Stratified on Geographic Region  Select Family Units | We could also use multistage sampling.<br>Here we view the population as having clusters… say family units are clusters in this example.  Each family is represented here by dots which are circled.<br>We select two family units from all family units.   The family units we select are marked in red.  Then within these family units (PSUs or clusters) we select all family members.  [ You could, alternatively, add another stage and select individuals within each family.] |
| Slide 27 | **Combine Different Sampling Designs**<br><br>• SRS, Systematic, Stratified and Multistage can be (usually are) combined<br><br>• Often helpful to diagram a study to understand the study design<br><br>UNC | As we have said a few times,  different sampling designs are often combined.<br>Sampling designs can be quite complicated, so often it is helpful to diagram a study design. |
| Slide 28 | **Quality of Life in Overweight Children in Australia**<br><br>Objectives: "To determine relationships between weight and health-related quality of life in a population sample of elementary school children."<br>Design: "Random stratified, 2-stage sampling design" 1569 children in Victoria, Australia were weighed and measured in 2000. Questionnaires were completed by child and parents assessing physical, emotional, social and school functioning.<br>Result: Negative effects on QOL with increasing child weight.<br>Williams, J., "Health Related Quality of Life of Overweight and Obese Children" JAMA,  January 5 2005, 293 -1 (pp. 70-76). | Let's turn to a global health example regarding the Quality of Life of Overweight Children in Australia which used many of the sampling topics we've discussed.<br>For the best understanding,  please read articles referenced at the end of the slide presentation if possible.  This is the same example that we've referred to in several of the lessons. The relevant portions are summarized for you, but you can read the articles listed in the reference for more information.<br><br>Very briefly, let's start with the objectives of the study.<br>The primary objective was to determine relationship between weight and health related quality of life in a population sample of elementary school children in Australia.<br>The design was a random two-stage sampling design.  1569 children in Victoria Australia were weighed and measured in the year 2000.  The children and the parents completed questionnaires about physical, emotional, social and school functioning.<br>The questionnaire results showed that increasing child weight was associated with decreased quality of life in this population.  More results will be given near the end of this presentation.<br>Let's look at the sampling design of this study and how it illustrates much of what we've learned (although I may make/change some assumptions about the original study or simplify the design to make a point about what we've talked about in sampling.) |

| Slide 29 |  STUDY DESIGN "A stratified two stage design was employed to select 24 primary schools. First schools were selected within each educational sector (Government, Catholic or Independent) with a probability proportional to size. Second, an entire class at each year level was randomly sampled within each school." | A quote from the article…..<br>"A stratified two-stage design was employed to select 24 primary schools. First schools were selected within each educational sector (Government, Catholic or Independent) with a probability proportional to size. Second, an entire class at each year level was randomly sampled within each school."<br>Let's see if we can get a mental picture of this description. Investigators obtained a sampling frame of all primary schools in Victoria.  The schools were categorized as either Government, Catholic or Independent.   These three categories were the strata.  Within each stratum, schools were selected for a total of 24 schools.  We cannot know from the information  that I've given so far how many schools were selected within each stratum – just that the total number of schools is 24. It could be that eight schools were selected in each stratum, or not.   In the diagram I show 8 school from each stratum, but we really don't know this from the article….. We also don't know how many schools were on the sampling frame. This is often the case in the real world – investigators are limited in the amount of information they can include in an article so the reader is left to wonder a bit.<br>So, we can see, so far, this is a stratified design.<br>Also, we can see that the study is a multistage (2 stage) design.  First schools are selected (they are the PSU's the primary sampling units).  The second stage is to select an entire class at each year level with in the  school.  So, classrooms are SSU (secondary sampling units).<br>The description indicates schools were selected with a "probability proportional to size".   This means that clusters (schools) which are bigger have a greater probability of being selected and that probability is proportional to the size of the cluster (school).<br>Some questions to think about so far….what are the sampling frames in this example?  Why did they use a multistage design?  Why did they use a stratified design? What are the implications if a school did not participate? Was that school replaced by another school in that stratum? Did they use SRS to select classrooms within the schools? How many primary schools were in Victoria and is 24 a pretty good representation?  (We will of course talk about sample size more formally, but I also think it is helpful to just "ball park" it -  24 school were selected out of how many?) How did they define 'primary school'?  In reading an article, or designing a study yourself, careful consideration of questions like this can be helpful (instructive?). |

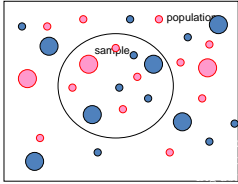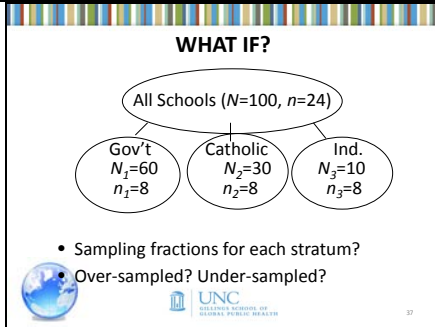| | | |
|---|---|---|
| Slide 30 | **Review: Select a SRS**<br><br>Suppose investigators randomly selected 5 independent schools from 72 independent schools. [They used probability proportional to size - bigger schools were more likely to be chosen]<br>Random number table<br>98663 89213 15388 74346 16125 30168 90229<br>One method: uniform random # * N =value take next largest number<br>1st school: 0.98663 * 72= 71.02 => 72<br>2nd school: 0.89213 * 72= 64.23 => 65<br>3rd school: 0.15388 * 72 = 11.08 => 12<br>4th school: 0.74346 * 72 = 53.53 => 54<br>5th school: 0.16125 * 72 = 11.61 => 12 duplicate<br>6th school: 0.30168 * 72 = 21.72 => 22<br>UNC KILLINGA SCHOOL OF GLOBAL PUBLIC HEALTH 30 | This is a review of how to select a SRS.<br>Let's suppose investigators randomly selected a SRS of independent schools from 72 independent schools.<br>[The actual study used probability proportional to size where bigger schools were more likely to be chosen- but we'll simplify that and use simple random sampling]<br>Let's list the first 5 schools that would be selected, out of the total sample size (like 8)<br>Once we got uniform random numbers from a random number table, we multiply the number by N and then take the next larger integer.<br>The first 5 schools to be selected are {72,65,12,54,22,…} |
| Slide 31 | **STATISTICS AND PARAMETERS**<br><br>"The achieved sample mirrored Victorian census data for age distribution, sex, ethnicity…"<br><br>• Percent female in sample → 49.5%<br>• Percent female in Victoria* pop'n → 50.9%<br>• Percent of overweight females in sample →21.3%<br>• Percent of overweight females in pop'n→ unknown<br><br>*from 2001 census  UNC KILLINGA SCHOOL OF GLOBAL PUBLIC HEALTH 31 | Let's talk a bit about statistics and parameters.<br>The article states: "The achieved sample mirrored Victorian census data for age distribution, sex, ethnicity…"<br>Why is this important? Why did the authors take the opportunity to tell the reader this? Well, the authors want to show that their sample is a good representation of the population of interest. In other words, from the census we know the age distribution, sex and ethnicity distribution etc."<br>If our sample is representative, we would expect the sample to have similar ages, genders and ethnic background as the population.<br>Suppose that:<br>The percent female in from the final sample was 49.5%.<br>And suppose that the percent female in Victoria was 50.9% (from 2001 census).<br>Here we are not using a statistic to estimate a parameter – the parameter is known. We are using the statistic(s) to show that the sample is pretty representative.<br>The percent of overweight females in sample is 21.3% according to the article.<br>The percent of overweight females in population is unknown.<br>[ We can talk about whether 50.9% and 49.5% are statistically significant, but for now let's just say that they appear to be pretty close.] |

| | | |
|---|---|---|
| Slide 32 |  | So to drive this point home, let's label them…<br>49.5 % female is a statistic (from the sample)<br>50.9 % female is a parameter (from the population)<br>21.3% overweight (within the females in the sample) is a statistic – it is measurable and describes the sample.. It varies from sample to sample…<br>And the percentage overweight in the population of females is unknown - this value is one of the values we wish to estimate.<br>Since the percent female in the sample is pretty close to the percent female in the population, we would like to think that the percent overweight in the sample is pretty close to the percent overweight in the population. It is a good sign that our sample is representative of the population. |
| Slide 33 |  | In conducting a sample, we want our sample to be representative of the population.<br>Suppose the population only had 30 students and we are take a SRS of 10 students out of 30 students.<br>In this picture, consider the pink dots to be female subjects, the blue dots to be males subjects, the big dots are overweight children and the small dots to be healthy weight children.<br>There are 15 boys and 15 girls in the population. Go ahead and count the dots. Make sure I'm right! So in this (small) population there are 50% girls and 50% boys.<br>If we've taken a representative sample then we'd expect our sample to have about the same proportions. Let's check.. In the sample (inside the ring), what is the percent of girls? (5 out of 10 or 50%) and the percentage of boys is 5 out of 10 or 50%. So the sample has the same gender distribution as the population. Our sample appears to be representative. Good!<br>Convince yourself that the proportion of overweight females in the sample is 1 out of 5 (20%). You'll probably need to pause the slide show at this point. In the sample, there are 5 females and one is overweight…. The statistic is 20% of females in the sample are overweight. In the population we can see that 3 out of 15 females are overweight (=20%).<br>Turning our attention to the fellas… 2 out of 5 in the sample are overweight (40%)… 6 out of 15 in the population are overweight (again 40%).<br>You can also convince yourself that the proportion overweight in the sample (3 out of 10) is the same as the proportion overweight in the population (9 out of 30).<br>This diagram is to show visually how a SRS should "look like" the population in the variable of interest. The sample should be representative of the population. |

| | | |
|---|---|---|
| Slide 34 | **SOURCES OF ERROR**<br><br>• 2336 students identified, 1943 completed the questionnaire. Suppose females were less likely to complete the questionnaire than males → *nonresponse error*<br>• <u>Suppose</u> the sampling frame omitted two schools within the Victoria area and included one school outside the Victoria → *frame error*<br><br>UNC<br>GILLINGS SCHOOL OF GLOBAL PUBLIC HEALTH<br>34 | Returning to the original study….<br>There are sources of error, let's look at a couple:<br>Out of 2336 students identified, 1943 completed the questionnaire. Suppose females were less likely to complete the questionnaire than males. What type of error may occur? – nonresponse error<br> Suppose the list of schools from which the sample was taken omitted two schools within the Victoria area and included one school outside the Victoria. This type of error is best described as what type of error?- frame error |
| Slide 35 | **STUDY DESIGN**<br><br>"A stratified two stage design was employed to select 24 primary schools. First schools were selected within each educational sector [Government, Catholic or independent] with a probability proportional to size. Second, an entire class at each level was randomly selected within each school."<br>What is 1st stage? – select a school<br>What is 2nd stage? – select class<br><br>UNC<br>GILLINGS SCHOOL OF GLOBAL PUBLIC HEALTH<br>35 | Let's go back to the study design<br>"A stratified two stage design was employed to select 24 primary schools. First schools were selected within each educational sector [Government, Catholic or independent] with a probability proportional to size. Second, an entire class at each level was randomly selected within each school."<br>Recall the 1st stage  -- select a school<br>The second stage is selecting a classroom within each grade level from the schools which we selected. |
| Slide 36 | **Why use multistage design?**<br>Eliminates need to list all students in Victoria<br>More efficient to administer questionnaire to entire classroom<br><br>**Why use stratified design?**<br>Assure selection within each subgroup<br>Adequate sample size for subgroup analysis<br>Allows for under- or over- sampling<br>Increase efficiency<br><br>UNC<br>GILLINGS SCHOOL OF GLOBAL PUBLIC HEALTH<br>36 | I asked these questions earlier,  now let's answer them.<br><br>Why do the investigators use multistage design?<br>        It eliminates need to list all students in Victoria.<br>        It is more efficient to administer questionnaire to entire classroom.<br>Why do the investigators use a stratified design?<br>        In order to assure selection within each subgroup.<br>        It can provide a adequate sample size for subgroup (school type) analysis.<br>        It allows for under or over sampling within the strata.<br>        We can increase efficiency if the strata are homogeneous.  You'll just have to trust me on that.<br><br>There are downsides, of course.   The analysis is more difficult.  Also with this clustered (multistage) design a bigger sample size is usually needed.  Again, we haven't gotten into the reasoning… just be aware there is "no free lunch"….<br>In statistics, as in most things, there are tradeoffs. |

| Slide 37 |  | Let's have some more fun with this example by making some assumptions…

Let's simplify this example in order to practice some of the material we've learned.
Let's just look at the first part of the design the stratification. We'll just consider this as a stratified design– with the strata being school type. The investigators pick schools from each stratum (Government, Catholic, Independent).

Further, for simplicity, let's suppose there are 100 schools in the ' population' of schools – it is nice round number. 100 primary schools in Victoria to choose from. We are going to select 24 schools out of 100 possible schools.
    Also, assume that there are 60 Government, 30 Catholic, 10 Independent (we don't actually know this number, but suppose this is the breakdown for conversation sake)
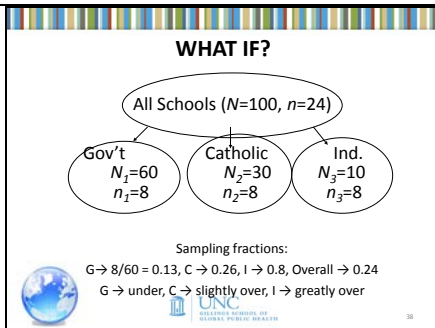    (You can think of the as proportions (0.6 Government, 0.3 catholic and 0.1 independent) or percentages (60% Gov't, 30% Catholic and 10% Independent).

Suppose we select 8 schools in each stratum. We've talked about previously how often investigators will make the sample size in each stratum be equal rather than the sampling fraction.
That brings up an interesting question… what are the sampling fractions?
Are there strata which are under or over sampled? Are we all comfortable with the notation so far? Try these questions on your own first and then I'll address them on the next slide. |
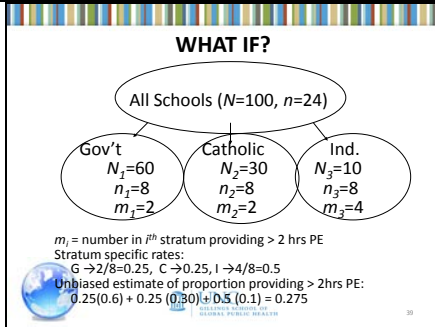| Slide 38 |  | Let's compute the stratum specific sampling fractions. We are going to calculate the sampling fraction within each group. …."stratum specific" means for each stratum (each group) sampling fraction (recall the sampling fraction is sample size/ population size or n/N).
For the gov't schools the sampling fraction is 8/60 = 0.13 (13%).
For the Catholic schools the sampling fraction is 8/30 or 26% and the independent schools have a sampling fraction of 0.8. The overall sampling fraction is 24/100 or 0.24.

The Government schools are undersampled (because their sampling fraction is less than 0.24. The independent schools are greatly oversampled -their sampling fraction is 0.8 much greater than 0.24. |

| Slide 39 |  | Further, suppose may we collect data at the school level. Let's review calculating an unbiased estimate for a proportion in a stratified sample.<br>We may ask them "Does your school provide greater than 2 hours of Physical Education for students per week?"<br>Let $m_i$ be the number of schools in the $i$th stratum in the sample that provide 2 hours of PE for student per week.<br>So in the government sample of 8 schools we might find that 2 schools provide this level PE (0.25)<br>In the Catholic schools in the sample, also 2 schools out of 8 provide > 2 hours of PE (0.25)<br>And among the 8 independent schools, 4 provide that level of PE (0.5)<br><br>Next we want to calculate an unbiased estimate of the proportion of schools in the population providing that level of PE. We know the proportion in the sample 8 out of 24. We want to know about the population of 100 schools…<br>Recall we weight stratum specific use of the program by the stratum specific sampling fractions. In English, we want to weight those values, 0.25 0.25 and 0.5 by the proportion in each stratum in the population.<br>$0.25(0.6) + 0.25(0.30) + 0.5(0.1) = 0.15 + 0.075 + 0.05 = 0.275$<br>What does this mean… in the population of 100 schools, we'd expect 0.275 of the schools to be providing > 2 hrs of PE… or approximately 28 schools out of 100. |

| | | |
|---|---|---|
| Slide 40 | **QOL of Overweight Kids in Australia**<br><br>• Survey (parents and child): PedsQL measured physical, emotional, social and school functions<br>• 20.2% overweight, 4.3% obese<br>• QOL Scores decreased (both parent- and child-reported) with increasing weight category for total score and subscores for physical, social, school and emotional functioning<br><br>UNC GILLINGS SCHOOL OF GLOBAL PUBLIC HEALTH<br>40 | So what did this actual stratified, 2 stage sampling design tell us about the Quality of Life and Obesity among Children in Australia?<br><br>Both the children and parents were given a survey called the PedsQL which measured physical, emotional, social and school functioning. The 1,456 children who participated in the survey were primarily 9- to 12- years old. Approximately 20.2% of children in the sample were classified as overweight, 4.3% were classified as obese. <u>Obesity is not just a US problem</u> but is an issue in other countries, particularly developed countries, as well.<br> The paper reports that "Parent and child perceptions were strikingly similar with obese children having the lowest summary and subscale scores." Quality of Life Scores decreased (both parent- and child- reported) with increasing weight category for total score, physical, social, school and emotional functioning (although not all decreases were statistically significant. )<br>Health related quality of life ( as measured by the PedsQL instrument) was lower for children in the heavier weight categories with the obese children in general having the lowest scores in each of the subcategories.<br>So childhood obesity (outside the US) is an important public health issue not only for the physical impact on the child's health but also for the emotional and psychological toll that associated with being overweight.<br>The sampling methods we've learned are used in many studies – such as this one with global impact – to select subjects in sample to tell us about associations in the population we are interested in.<br>[Williams, J., "Health Related Quality of Life of Overweight and Obese Children" JAMA, January 5 2005, 293 -1 (pp. 70-76).] |
| Slide 41 | **Objectives**<br><br>▪ Review Simple Random Sampling, Stratified Sampling and Multi-stage Sampling<br>▪ Understand the role of sampling in motivating example: Quality of Life and Obesity in School Children in Australia<br><br>UNC GILLINGS SCHOOL OF GLOBAL PUBLIC HEALTH | |

**Slide 42**

## REFERENCES

Introduction to the Practice of Statistics, 6th edition, Moore and McCabe, W.H. Freeman and Company.

Survey Sampling, Kish, John Wiley and Sons publishing, 1995.

Williams, J., "Health Related Quality of Life of Overweight and Obese Children" JAMA, January 5 2005, 293 -1 (pp. 70-76).

Waters, E., "The Child Health Questionnaire in Australia: reliability, validity and population means" Australian and New Zealand journal of Public Health, April 2000, 24-2 (pp. 207-210).

UNC
GILLINGS SCHOOL OF
GLOBAL PUBLIC HEALTH

42