

## Unit 5 Populations and Samples

*“To all the ladies present and some of those absent”*  
- Jerzy Neyman

The collection of all individuals with HIV infection and the collection of all individuals with exposure to mercury are examples of **populations**. A **census** involves the collection of information on **every individual** in the population and is one way to obtain information about a population. How nice! Precisely because we have the information needed for every individual in the population, there is no need for any statistics at all.

Unfortunately, most of the time, censuses are impractical because we lack the necessary resources to obtain the sought for information on every individual.

So, instead of the (ideal) census, we study a subset of the population, called a **sample**. We can calculate numbers from a sample (these are called statistics) and these are used to make (hopefully meaningful) inferences about the population.

There are lots of ways to obtain a sample of a population; these are called **sampling designs**. Perhaps the most familiar is the method of **simple random sampling**. Loosely, simple random sampling is sampling at random without replacement from the population.

The goal of a **sampling design** and the statistical analyses that follow are: (1) to obtain a sample with a known probability of selection and for which the conclusions drawn are (2) in the long run correct (unbiased) and (3) in the short run in error by as little as possible (minimum variance).

Cheers!

## Dragon Sampling II



The IRB won't be happy with our use of involuntary sampling

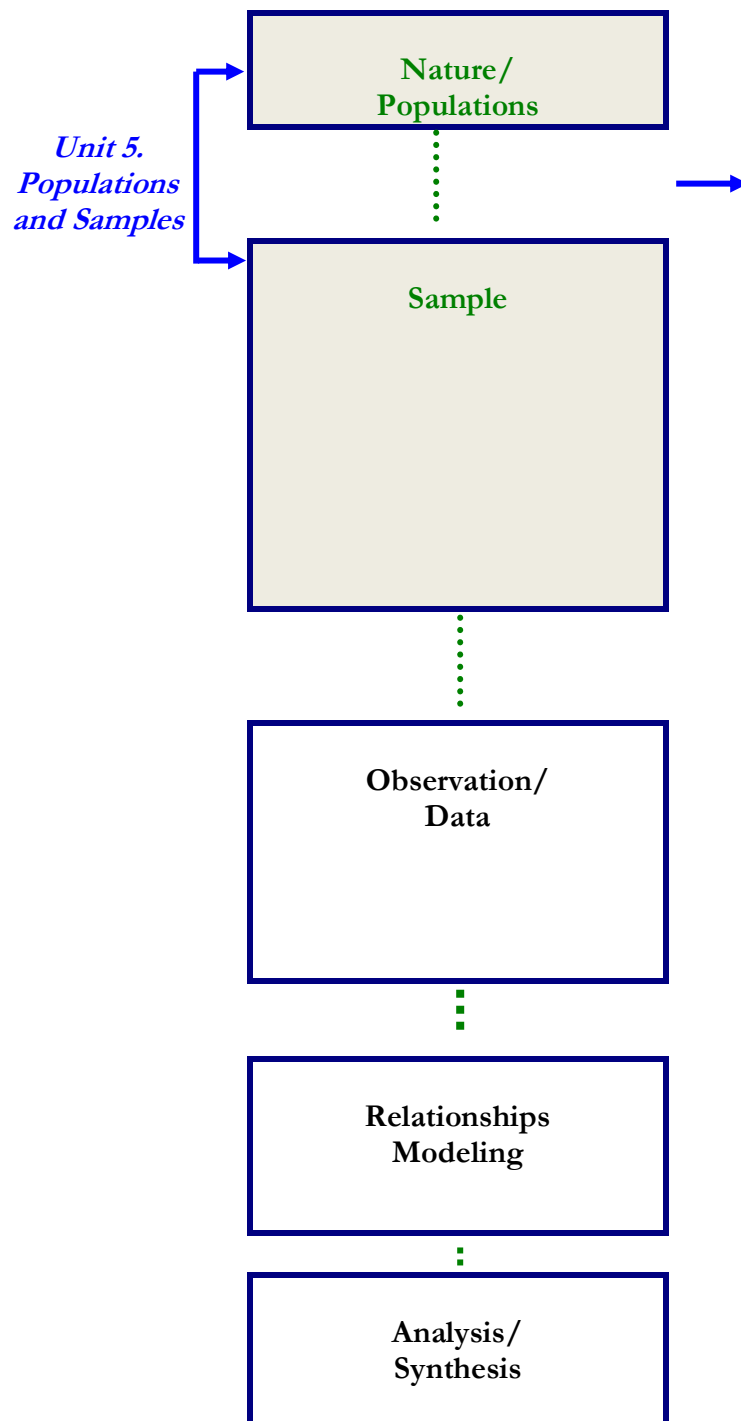
(Source: J.B. Landers. With permission, download from CAUSEweb.org)

Nature ——— Population/  
Sample ——— Observation/  
Data ——— Relationships/  
Modeling ——— Analysis/  
Synthesis

## Table of Contents

<b>Topics</b>	1. Unit Roadmap .....	4
	2. Learning Objectives .....	5
	3. A Feeling for Populations v Samples .....	6
	4. Target Populations, Sampled Populations, Sampling Frames .....	9
	5. On Making Inferences from a Sample .....	13
	6. Simple Random Sampling .....	16
	7. Some Non-Probability Sampling Plans .....	19
	8. Some Other Probability Sampling Plans .....	22
	a. Systematic .....	22
	b. Stratified .....	24
	c. Multi-stage .....	26
	9. The Nationwide Inpatient Survey (NIS) .....	27
<b>Appendix</b>	More on Simple Random Sampling .....	28
	a. Sampling WITH v WITHOUT replacement .....	31
	b. How to select a simple random sample WITHOUT replacement..	38

## 1. Unit Roadmap



**Take another look at the roadmap at the footer of this page.** A study begins with a sample from a population (highlighted here in bold red). → From our sample, we make some observations and record some data. → From there, with the use of an assumed model, we estimate some things (for example – average response to treatment in the population) and test some hypotheses (for example – the new treatment is no better than the control treatment).

In the real world, we have just the *one sample* and no luxury to repeat the study over and over. So, we rely on the properties of the sampling procedures (for example – all theoretically possible samples were equally likely to have been selected) as the justification for the conclusions we draw.

**Unbiased** - If we were to repeat our study over and over again, the average of our sample (for example – the average response to treatment) will eventually settle down to a long range average of the average that is equal to the true population average.

**Minimum Variance** –A conclusion drawn from a sample will differ from the reality of the population. This is *sampling error*. An additional goal of sampling is to obtain a sample for which sampling error is minimized.

## 2. Learning Objectives

When you have finished this unit, you should be able to:

- Explain the distinction between target population, sampled population, and sample.
- Explain why it is important to strive for a sample that is representative of the population from which it is taken.
- Explain the rationale for choosing a sampling method that minimizes sampling error.
- Distinguish non-probability versus probability samples.
- Define simple random sampling.
- Explain the rationale for systematic, stratified, and multi-stage sampling methods.
- Define systematic, stratified, and multi-stage sampling.

Interested readers, reading the appendix, will also have a feel for:

- The distinction between sampling with versus without replacement.

### 3. A Feeling for Populations versus Samples

In [Units 1 and 2](#), our goal was to summarize (and communicate effectively!) the information in a given sample. We didn't concern ourselves with the source of the sample. That is, at that point, we didn't give any thought to the population from which the sample was obtained. In Units 1 and 2, we learned about various kinds of summaries (graphical and numerical).

In [Units 3 and 4](#), we began to think about populations and random draws from populations. Unit 3 was an introduction to the most basic of probability distributions and we considered the simple scenario of all random draws of elementary outcomes being equally likely. We learned about events, that they can be either elementary outcomes or collections of elementary outcomes. We also learned some tools of basic probability calculations, depending on the type of events: mutually exclusive, dependent/conditional, independent, etc. In Unit 4, we revisited some common epidemiologic study designs (and their analysis approaches) as examples of conditional probability distributions (for example the cohort design calls for sampling from each of two independent distributions, one set from the distribution of exposed and the other set from the distribution of non-exposed).

Here in [Unit 5](#), we put the two together: **population** and **sample**. A sample is obtained as the result of following a sampling procedure. The nature and specifics of the sampling procedure is called the **sampling design**.

**Meaningful statistical inference requires that the sample studied be a probability sample.**

- Population – The collection of all the individuals of interest.
- Probability Sampling Design – The rules of probability that govern the likelihood of each sample being selected
- Sample – The subset of the population that is selected as the result of sampling.

Nature ——— Population/  
Sample ——— Observation/  
Data ——— Relationships/  
Modeling ——— Analysis/  
Synthesis

Non-representative sampling may (and often does)  
produce study conclusions that are incorrect.

### Example – The 2016 Presidential Election

*Let's not talk about the 2020 election*

- Before the 2016 presidential election, there were numerous polls of “anticipated” voters.
- Each was asked whom they were going to vote for – Clinton or Trump

	Predicted by 538 Chance of Winning	True Election Result # Electoral Votes
Clinton	71.4%	227
Trump	28.6%	304

- How could the prediction have been so wrong?
- How could the samples have been so dissimilar to the voting population?

Nature ——— Population/  
Sample ——— Observation/  
Data ——— Relationships/  
Modeling ——— Analysis/  
Synthesis

## The 2016 Presidential Election

- The outcome of the election is determined by the Electoral College, not the popular vote. *(yes, yes ... you knew that already...)*
- Polls are a “snapshot” of a point in time. As such, they are not necessarily a good predictor of future behavior. *(including such things as: 1) changing one’s mind about voting or not voting; and 2) changing one’s preference from Clinton to Trump or vice versa...)*
- Predictions from polls are also based on models and these may incorporate assumptions that are incorrect. *(E.g. – the profiles of who voted in past elections were not representative of who chose to vote in the 2016 election; that is “Assumed model of likely voters  $\neq$  Actual voters” ...)*

Thus, bias occurred because:

- (1) there was over-sampling *(and probably some under-sampling, too); and*
- (2) the nature of the over *(and/or under-sampling)* was related to voter preference.

**IMPORTANT POINT** – “Oversampling”, per se, does not produce bias in study findings necessarily. For bias to occur, the disproportionate sampling has to be related to the outcome. For example, consider “oversampling” people whose favorite ice cream is chocolate chip.... Note – We will say much more about “bias” later. For now, think of bias as the extent to which a finding is incorrect.

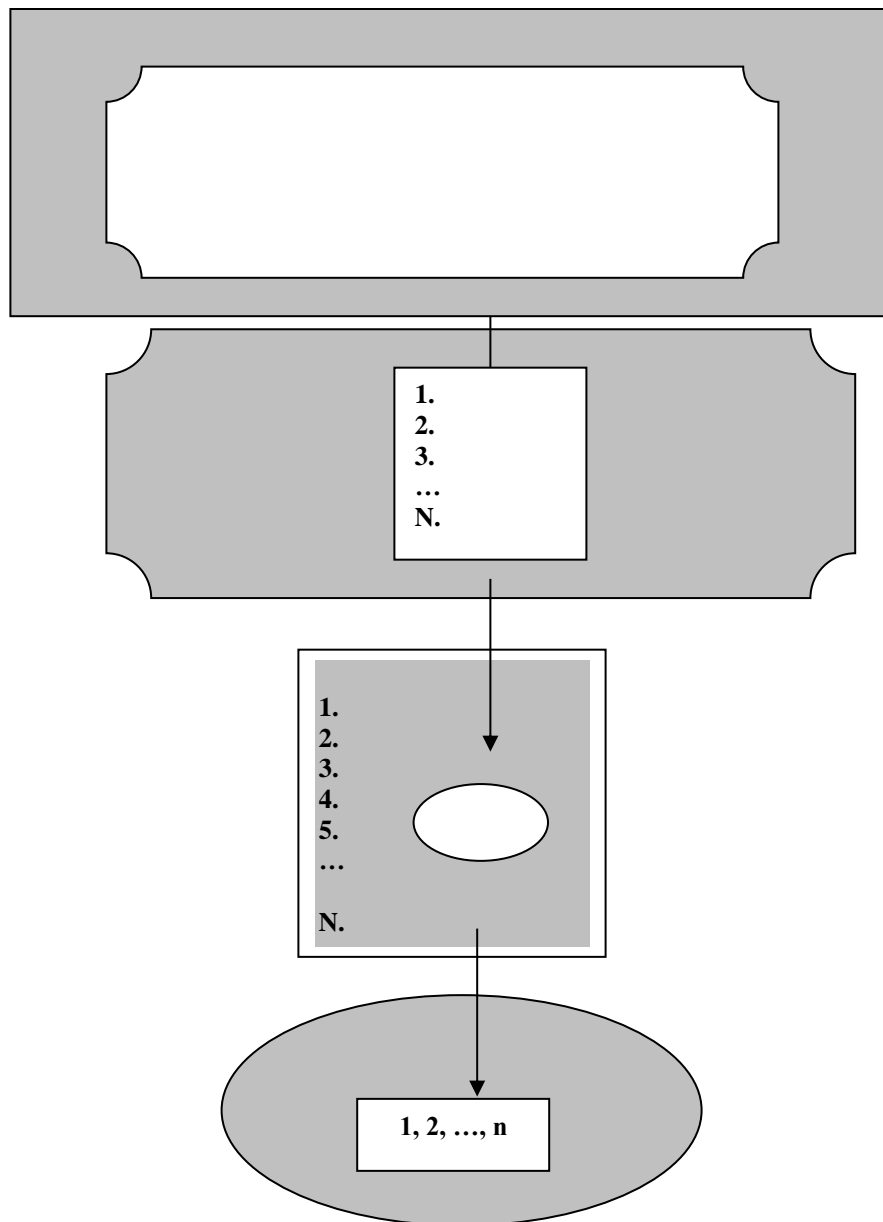
### Example, continued - The 2016 Presidential Election

The populations actually sampled (the sampled populations) were *not the same* as the population of interest (the target population comprised of those individuals who actually cast a vote)

Nature ——— Population/  
Sample ——— Observation/  
Data ——— Relationships/  
Modeling ——— Analysis/  
Synthesis



#### 4. Target Population, Sampled Population, Sampling Frame



##### Target Population

The whole group of interest.

*Note – A convention is to use capital “N” to represent the size of a finite population.*

##### Sampled Population

The subset of the target population that has at least some chance of being sampled.

##### Sampling Frame

An enumeration (roster) of the sampled population. **So, yes. The sampling frame and the sampled population are the same thing. The one that has been put into a list is the sampling frame.**

##### Sample

The individuals who were actually measured and comprise the available data.

*Note – A convention is to use small “n” to represent the size of a sample.*

Target Population	<ul style="list-style-type: none"> <li>• The entire collection of individuals who are <b>of interest</b>.</li> <li>• <b>Example</b> – The population of 2016 presidential election registered voters who <i>actually</i> voted.</li> </ul>
Sampled Population	<ul style="list-style-type: none"> <li>• The aggregate of individuals that was <b>actually sampled</b>.</li> <li>• A listing of the entire sampled population comprises the <b><u>sampling frame</u></b>.</li> <li>• <b>GOAL:</b> sampled = target</li> <li>• The sampled population is often difficult to identify. We need to ask: Whom did we miss?</li> <li>• Constructing the sampling frame can be difficult</li> </ul>

## Sampling Frames Why They Are Difficult

### To Construct a Sampling Frame Requires

- ♣ Enumeration of every individual in the sampled population
- ♣ Attaching an identifier to each individual
- ♣ (Often, this identifier is simply the individual's position on the list)

### Example –

- ♣ The League of Women's Voters Registration List might be the sampling frame for the target population who voted in the 2020 presidential elections.
- ♣ Individual identification might be the position on this list.

### Now You Try –

- ♣ The target population is joggers aged 40-65 years.
- ♣ How might you define a sampling frame?

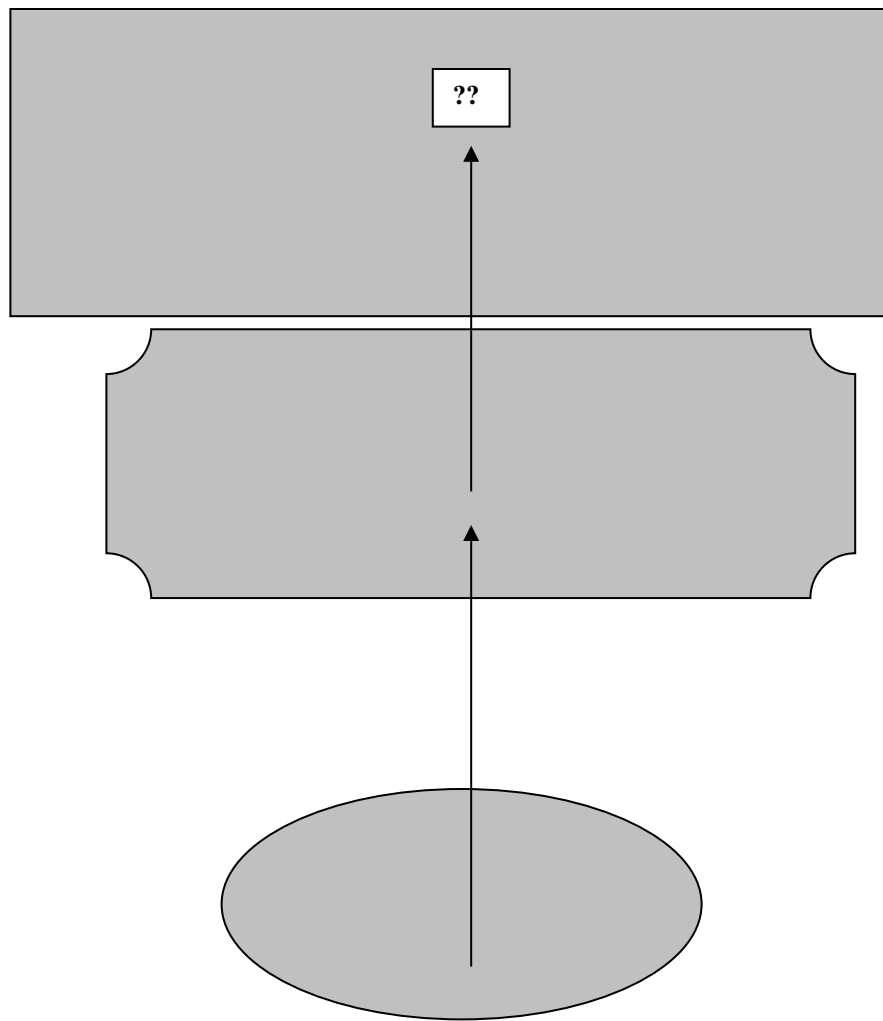
## HOMEWORK DUE Friday October 21, 2022

### Question #1 of 2

For each of the following situations, define the target population, and how you might obtain a sample. What will be your sampled population? How does this differ from the target population?

- a. A city engineer wants to estimate the average weekly water consumption for single-family dwelling units in the city.
- b. A physician wants to estimate the average length of time from initial diagnosis with ovarian cancer to death.

## 5. On Making Inferences From a Sample (this time - read from the bottom up)



### Target Population

The conclusion may or may not generalize to the target population.

- Refusals
- Selection biases

### Sampled Population

If sampling is representative, then the conclusion generalizes to the sampled population.

### Sample

The conclusion is drawn from the sample.

- ♣ The conclusion is initially drawn from the sample.
- ♣ The question is then: How far back does the generalization go?
- ♣ The conclusion usually applies to the sampled population
- ♣ It may or may not apply to the target population
- ♣ The problem is: It is not always easy to define the sampled population

### **Example – NIH Funded Randomized Trial**

- ◆ The sampling frame, by definition, is allowed to contain only consenters
- ◆ Thus, refusers, by definition, are not in the sampling frame.
- ◆ Thus, in a randomized trial protocol that includes consent, the sampled population differs from the target population because the sampled population is restricted to consenters only.
- ◆ This suggests that in any study, the preliminary analyses should always include a comparison of the consenters versus the refusers.

### ***Now You Try ...***

- ◆ Suppose the target population is current smokers.
- ◆ How might you construct a sampling frame?
- ◆ What do you end up with for a sampled population?
- ◆ Comment on the nature of generalization, to the extent possible.

## HOMEWORK DUE Friday October 21, 2022

### Question #2 of 2

Which of the following estimates are likely to be biased? Why? Is the bias positive or negative? Why? (*note: Positive bias means a consistent likelihood of overestimating, negative bias is underestimating*).

- You estimate the average number of bank customers waiting for service whenever the bank is open by counting the number of customers whenever you go to the bank.
- You estimate the proportion of 7-12 year old children using helmets when they ride bikes by asking parents if their child wears a helmet when the child is brought to the physician's office for a "well" visit.
- A highway patrolman parks next to a highway and records speeds on his radar to estimate the percentage of people exceeding the speed limit on that highway.

## 6. Simple Random Sampling

We would like our sampling plan to produce estimates that are:

- ◆ **Unbiased** – If sampling is repeated over and over and over, the “long run” average conclusion about the population should be correct.
- ◆ **Minimum variance** – The discrepancies (and their variance) between the conclusions drawn from the separate samples versus what is true in the population should be as small as possible;

Definition simple random sampling:

**Simple random sampling** is the method of sampling in which every individual in the sampling frame has the same (equal) chance of being included in the sample.

The virtue of simple random sampling is that it is **unbiased**, meaning:

- ◆ IF we draw sample after sample after sample after sample ....  
AND IF, for each sample, we compute a sample  $\bar{X}$  as our guess of  $\mu$ ,  
so as to compile a collection (yes - “**sampling distribution**”) of all possible sample estimates  $\bar{X}$ ,
- ◆ THEN “in the long run”...  
the average of all the sample estimates, average of ( $\bar{X}$  after  $\bar{X}$  after  $\bar{X}$  ...) will be equal to the population parameter value (the true value of  $\mu$ )



**Example of Simple Random Sampling, without replacement**  
*“Simple Random Sampling Without Replacement is unbiased”*

Suppose the Sampling Frame == Target Population

Subject ID	1	2	3	4	5	6
Age, years	21	22	24	26	27	36

In this (admittedly tiny) illustration, we can actually calculate the value of the population mean parameter. This is handy, as we’ll refer back to it later for the purposes of illustrating unbiased:

- ♦ Population mean age  $m = \frac{21 + 22 + 24 + 26 + 27 + 36}{N = 6} = 26$  years
- ♦ **Remember:** The investigator doesn’t know this value. That’s why they’re taking a sample!

The Investigator Executes the Following Sampling Procedure

- ♦ Draw a random sample of  $n=3$  subjects from the population. Do each successive draw “without replacement”. **Note** – “Without replacement” means that each selected person, once selected, is NOT returned to the population for future sampling. More on this later.

Calculate the sample mean, over and over again, once for each possible sample:

- ♦ For each sample, sample mean  $\bar{X} = \frac{\text{1st value} + \text{2nd value} + \text{3rd value}}{n = 3}$
- ♦ In this illustration (but not in real life because, remember, in real life the investigator does not know  $\mu$ ) we can calculate the error of each  $\bar{X}$  as an estimate of  $\mu$  by computing

$$\text{error} = (m - \bar{X}) = (26 - \bar{X})$$

How many different samples of size  $n=3$  are possible from a population of size  $N=6$ , when the sampling design calls for simple random sampling without replacement? To understand what it meant by a sampling design being unbiased, we want to consider all these samples. There are 20 possible “drawn without replacement” samples of size 3 from a population of size 6.

The table below shows all 20 samples. For each, it shows the three sampled values of age (2<sup>nd</sup> column), the sample average age (3<sup>rd</sup> column) and the discrepancy (“error”) between the sample mean and the population mean (3<sup>rd</sup> column):

Sample #	Sample (each of $n=3$ )	Sample mean, $\bar{X}$	Error = $\mu - \bar{X} = (26 - \bar{X})$
1	{ 21, 22, 24 }	22.333	+3.667
2	{ 21, 22, 26 }	23	+3
3	{ 21, 22, 27 }	23.333	+2.667
4	{ 21, 22, 36 }	26.333	-0.333
5	{ 21, 24, 26 }	23.667	+2.333
6	{ 21, 24, 27 }	24	+2
7	{ 21, 24, 36 }	27	-1
8	{ 21, 26, 27 }	24.667	+1.333
9	{ 21, 26, 36 }	27.667	-1.667
10	{ 21, 27, 36 }	28	-2
11	{ 22, 24, 26 }	24	+2
12	{ 22, 24, 27 }	24.333	+1.667
13	{ 22, 24, 36 }	27.333	-1.333
14	{ 22, 26, 27 }	25	+1
15	{ 22, 26, 36 }	28	-2
16	{ 22, 27, 36 }	28.333	-2.333
17	{ 24, 26, 27 }	25.667	+0.333
18	{ 24, 26, 36 }	28.667	-2.667
19	{ 24, 27, 36 }	29	-3
20	{ 26, 27, 36 }	29.667	-3.667
		Average = 26	Sum = 0

We have two ways of seeing that this sampling design is **unbiased**: 1) the average of the sample means is equal to the population mean; and 2) the errors “balance out” to zero.

(1) The average of the sample averages  $\bar{\bar{X}}$ , taken over all 20 possible samples, is  $\mu = 26$ .

(2) The sum of the errors,  $(\mu - \bar{X}) = (26 - \bar{X})$ , is 0.

$$\frac{\sum_{\text{all 20 possible samples}} \bar{X}}{20} = \mu = 26$$

$$\sum_{\text{sample \#1}}^{\text{sample \#20}} [\text{error}] = 0$$

Nature ——— Population/  
Sample ——— Observation/  
Data ——— Relationships/  
Modeling ——— Analysis/  
Synthesis

## 7. Some Non-Probability Sampling Plans

**Non-probability samples** are not random draws. Consequently, analyses of non-probability samples cannot be assumed to be representative of the population of interest. Even so, non-probability sampling methods are sometimes used (e.g., as when you get an email survey after visiting the dentist). Here are three examples of non-probability sampling plans:

- (1) Quota
- (2) Judgment
- (3) Volunteer/Convenience

### (1) Quota Sampling Plan

#### Example –

Population is 10% African American

Sample size of 100 must include 10 African Americans

#### How to Construct a Quota Sample

1. Determine the relative frequencies of each characteristic (e.g. gender, race/ethnicity, etc) that is hypothesized to influence the outcome of interest.
2. Select a fixed number of subjects of each characteristic ( e.g. males or African Americans) so that

Relative frequency of  
characteristic in **sample**  
(e.g. 10%)

MATCHES

Relative frequency of  
characteristic in **population**  
(e.g. 10%)

Nature ——— Population/  
Sample ——— Observation/  
Data ——— Relationships/  
Modeling ——— Analysis/  
Synthesis

## (2) Judgment Sampling Plan

Decisions regarding inclusion or non-inclusion are left entirely to the investigator. Judgment sampling is sometimes used in conjunction with quota sampling.

### Example

“Interview 10 persons aged 20-29, 10 persons aged 30-39, etc”  
Sample size of 100 must include 10 African Americans.

### Example

Market research at shopping centers

## (3) Volunteer/Convenience Sampling Plan

Volunteers are recruited for inclusion in the study by word of mouth, sometimes with an incentive of some sort (eg. gift certificate at a local supermarket)

### Example

For a study of a new diet/exercise regime, volunteers are recruited through advertising at local clinics, health clubs, media, etc.

*Problem?*

Nature ——— Population/  
Sample ——— Observation/  
Data ——— Relationships/  
Modeling ——— Analysis/  
Synthesis

### Limitations of a Non-Probability Sampling Plan

*They're serious!*

1. We have no idea if the sampling plan produces unbiased estimators. It probably doesn't.
2. Any particular sample may be highly unrepresentative of the target population.
3. Statistical inference, by definition based on some sort of probability model, is not possible.
4. Regarding quota sampling:
  - We have no real knowledge of how subjects were selected (recall the 2016 Presidential election).
5. Regarding judgment sampling:
  - Likely, there is bias in the selection – at least for the reasons of comfort and convenience.
6. Regarding volunteer sampling:
  - Volunteers are likely to be a “select” group for being motivated (or cranky and wanting to express it).

## 8. Some Other Probability Sampling Plans

### a. Systematic Sampling

- Population size =  $N$ . Desired sample size =  $n$ .
- The percent of the population we want in our sample is thus  $p = (n/N)$ .  
Thus, we want a  $(p)(100)$  % sample
- Step 1: Pick the first item by simple random sampling.  
Steps 2 onward: Thereafter, select every  $(N/n)^{\text{th}}$  item

#### Example:

Suppose  $n=20$  is desired from a population of size  $N=100$ .

→ This is a 20% sample, since  $p = (n/N) = (20/100) = .20$  or 20%.

Step 1: Pick the first individual by simple random sampling

Steps 2 onward: Thereafter, select every  $(N/n)^{\text{th}} = (100/20) = 5^{\text{th}}$  individual by systematic sampling. The first individual is selected by simple random sampling (so chances of inclusion are 1 in 100) Thereafter, take every 5<sup>th</sup> individual (so chances are then 0 or 1 depending on position in list!!)

#### Example –

Suppose we want a sample of size  $n=100$  from the  $N=1000$  medical charts in a clinic office.

→ This is a 10% sample, since  $p = (n/N) = (100/1000) = .10$  or 10%.

Step 1: Pick the first individual by simple random sampling

Steps 2 onward: Thereafter, select every  $(N/n)^{\text{th}} = (1000/100) = 10^{\text{th}}$  individual by systematic sampling. That is, thereafter, select every 10th chart.

### Remarks on Systematic Sampling

#### Advantages:

- It's easy.
- Depending on the listing, the sampled items are more evenly distributed.
- As long as there is no association with the order of the listing and the characteristic under study, this should yield a representative sample.

**Disadvantage:**

- If the sampling frame has periodicities (a regular pattern) and the rule for systematic sampling happens to coincide, the resulting sample may not be representative.

**Example of a Periodicity that Results in a Biased Sample:**

- Clinic scheduling sets up 15 minute appointments with physicians
- Leaves time for an emergency, or walk-in visit at 15 minutes before the hour, every hour.
  - Doing a chart-audit, you sample every 4th visit and get only the emergency visits selected into the sample, or else none of them.

**b. Stratified Sampling**  
*Simple Random Sampling within Strata*  
 (think subgroups)

Example-

Do construction workers experience major health problems?

Do health problems differ among males and females?

Construction workers, as a group, are likely to be comprised predominately of males.

Thus, if we take a simple random sample we may get very few women in the sample.

Procedure:

1. Define mutually exclusive strata such that the outcome of interest is likely to be

similar within a stratum; and very different between strata.

outcome:	health problems
strata:	Males at birth / Females at birth

2. Obtain a simple random sample from each stratum

We want to be sure to get a good overall sample.  
 Sampling each stratum separately ensures this.



## Remarks on Stratified Sampling

Advantage:

- Good when population has high variability, especially when the population includes a mix of people (e.g. males and females) that are NOT similarly represented (eg. population is disproportionately male)

Take care:

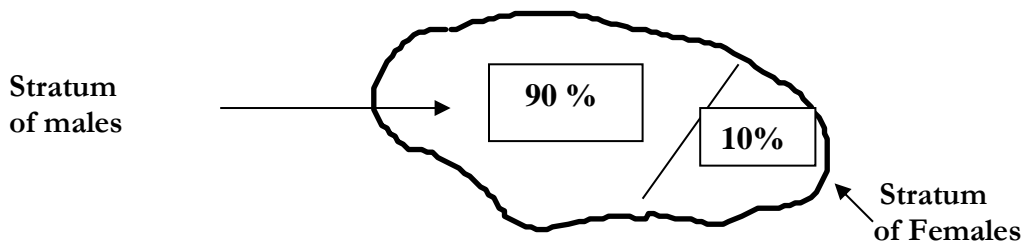
- The strata MUST be mutually exclusive (no overlap) and exhaustive (cover all possibilities)
- To compute an overall population estimate requires use of weights that correspond to representation in the population. The following is an example.

### Example of Calculation of Weighted Mean from a Stratified Sample -

**Goal** – Suppose we’re told that the average # cigarettes smoked per day is **38** (just shy of 2 packs!) among male construction workers and **12** (a little over a half a pack) among female construction workers. The goal is to estimate the average # cigarettes smoked per day among all construction workers.

The estimated mean should take into account the fact that the population is disproportionately male (90% male, 10% female)

- ♣ Since males are 90% of the population, let's weight the average for males using  $w = 0.90$
- ♣ Since females are 10% of the population, let's weight the average for females using  $w = 0.10$
- ♣ Note that weights total 1.00



$$\begin{aligned} \left[ \begin{array}{c} \text{Weighted} \\ \text{average, } \bar{X}_w \end{array} \right] &= \left( \begin{array}{c} \text{weight} \\ \text{males} \end{array} \right) (\bar{X}_{\text{males}}) + \left( \begin{array}{c} \text{weight} \\ \text{females} \end{array} \right) (\bar{X}_{\text{females}}) \\ &= (.90)[38 \text{ cigarettes}] + (.10)[12 \text{ cigarettes}] \\ &= 35.4 \text{ cigarettes/day on average, overall.} \end{aligned}$$

### c. Multi-Stage Sampling

*Good, Sometimes Essential, for “Difficult” Populations*

**Example** Suppose we want to study a gypsy moth infestation.

A multistage sample plan calls for

- 1ST - Select individual trees  
(Primary sampling units - PSU's)
- 2nd - Select leaves from only the selected trees  
(Secondary sampling units)

### Multistage Sampling

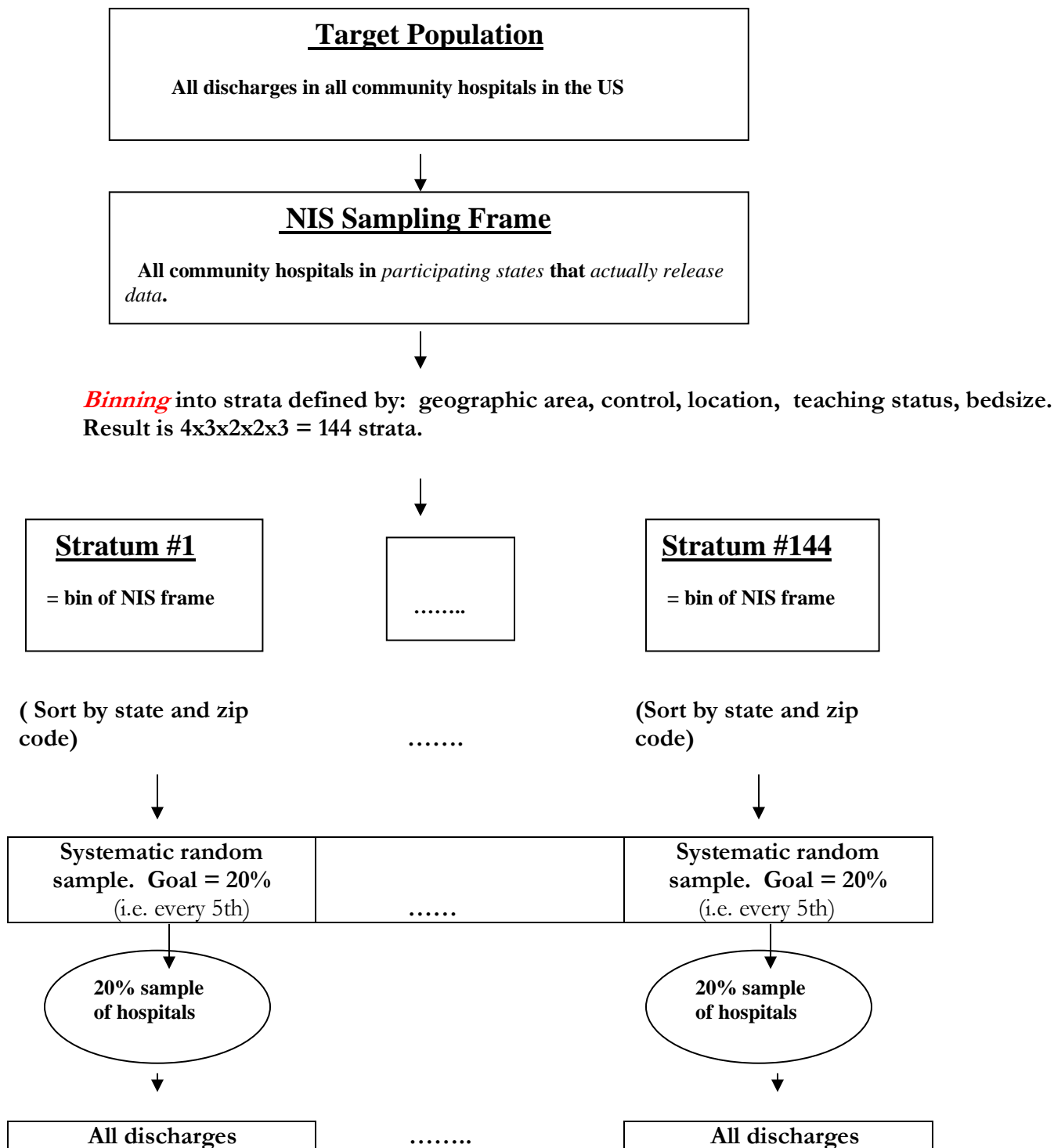
- The selection of the primary units may be by simple random sampling
- The selection of the secondary units may also be by simple random sampling
- Inference then applies to the entire population

### CAUTION!!!

- Take care that the selection of primary sampling units is NOT on the basis of study outcome. Bias would result.

Nature ——— Population/  
Sample ——— Observation/  
Data ——— Relationships/  
Modeling ——— Analysis/  
Synthesis

# 9. The Nationwide Inpatient Survey (NIS) Sampling Designs Can Be Quite Complex



## Appendix More on Simple Random Sampling

A **probability sampling plan** is “*out of the hands*” of the investigator.

- ◆ Each individual has a **known probability** of inclusion in the sample, prior to sampling.
- ◆ The investigator has **no discretion** regarding the inclusion or exclusion of an individual
- ◆ This eliminates one source of potential bias – that on the part of the investigator.

**How do we know if a sample is REPRESENTATIVE?**

- ◆ Ultimately, we don’t know!.
- ◆ So, instead, we use an unbiased sampling plan and hope for the best.
- ◆ In the meantime, we can generate some descriptive statistics and compare these to what we know about the population.

**Simple Random Sampling** is the basic probability sampling method.

Recall its definition from page 13:

**Simple random sampling** is the method of sampling in which every individual in the sampling frame has the same (equal) chance of being included in the sample.

$$\text{Under simple random sampling} \\ \text{Probability \{each equally likely sample\}} = \frac{1}{\text{number of equally likely samples}}$$

Thus, we need to solve for **the number of equally likely samples!**

**There are two kinds of sampling: WITHOUT replacement and WITH replacement**

- ◆ **Example of selection without replacement** – You are selected to participate in a survey. You can only be surveyed once .
- ◆ **Example of selection with replacement** – You play the lottery multiple times and so you are available for selection multiple times (not sure this is a great example ...).

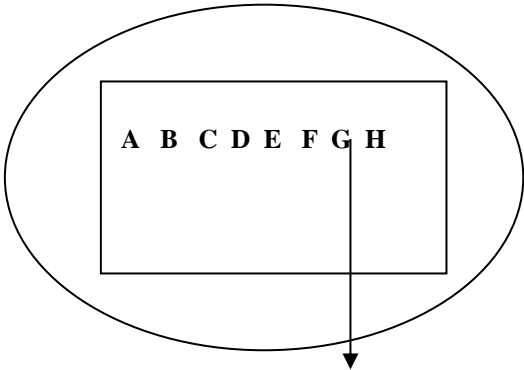
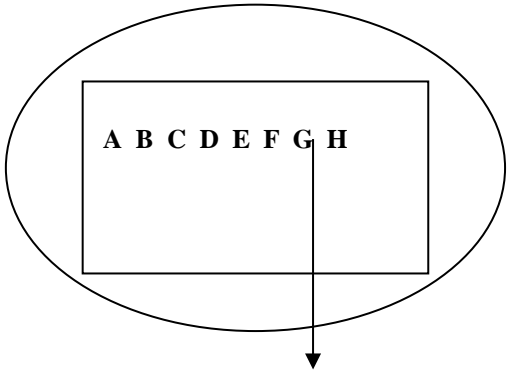
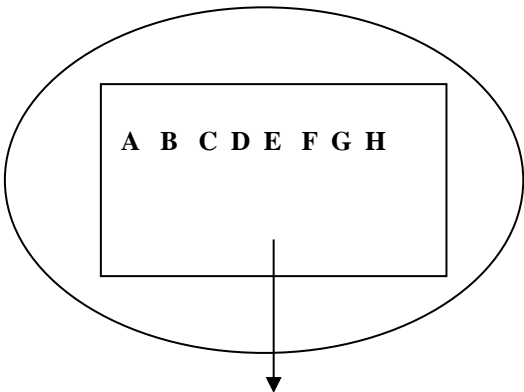
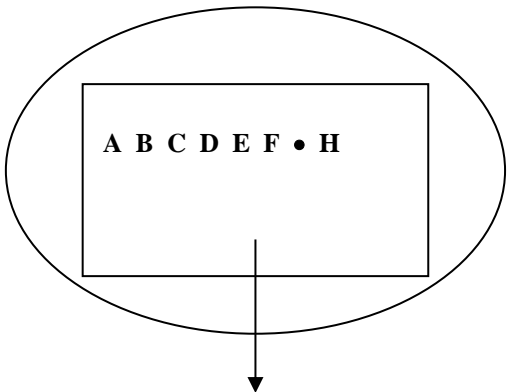
Nature ——— Population/  
Sample ——— Observation/  
Data ——— Relationships/  
Modeling ——— Analysis/  
Synthesis

**How many Equally Likely Samples Are there?**

a. Simple Random Sampling **With** Replacement

*versus*

Simple Random Sampling **Without** Replacement

With Replacement Answer: $N^n$	Without Replacement Answer: $N (N-1) (N-2) \dots (N-n+1)$
	
1 <sup>st</sup> Draw: A B C D E F G H are available - Suppose “G” is selected for inclusion	1 <sup>st</sup> Draw: A B C D E F G H are available - Suppose “G” is selected for inclusion
	
2 <sup>nd</sup> Draw: A B C D E F G H are ALL available - Thus “G” is available for inclusion a 2 <sup>nd</sup> time.	2 <sup>nd</sup> Draw: A B C D E F H are available but “G” is not available anymore - Thus, “G” can only be included once.
Etc	etc



a. Simple Random Sampling With versus Without Replacement  
General

With Replacement	Without Replacement
Population size = N Sample size = n	Population size = N Sample size = n
Each individual has a 1/N chance of inclusion in the sample.	Each individual has a 1/N chance of inclusion in the sample, <u>overall</u> .
How many equally likely samples are there? Answer: $N^n$	How many equally likely samples are there? Answer: $(N)(N-1) \dots (N-n+1)$
$\begin{array}{ccccccc} (N) & & (N) & & (N) & \dots & (N) \\ \uparrow & & \uparrow & & \uparrow & & \uparrow \\ 1^{\text{st}} & & 2^{\text{nd}} & & 3^{\text{rd}} & & n^{\text{th}} \\ \text{draw} & & \text{draw} & & \text{draw} & & \text{draw} \end{array} = N^n$	$\begin{array}{ccccccc} (N) & & (N-1) & & (N-2) & \dots & (N-n+1) \\ \uparrow & & \uparrow & & \uparrow & & \uparrow \\ 1^{\text{st}} & & 2^{\text{nd}} & & 3^{\text{rd}} & & n^{\text{th}} \\ \text{draw} & & \text{draw} & & \text{draw} & & \text{draw} \end{array} = (N)(N-1) \dots (N-n+1)$
Probability {each equally likely sample of size n} = $\frac{1}{N^n}$	Probability {each equally likely sample of size n} = $\frac{1}{N(N-1)(N-2) \dots (N-n+1)}$

Again, under simple random sampling:

$$\text{Probability \{each equally likely sample\}} = \frac{1}{\text{number of equally likely samples}}$$

**Example – How many equally likely samples are there?**  
**Simple Random Sampling WITH Replacement**

**Population**

Four Queens in a deck of cards

**Sampling Plan**

- Draw one card at random
- Note its suit
- *Return the selected card (“with replacement”)*
- Draw one card at random
- Note its suit

**Population size,  $N=4$**

**Sample size,  $n=2$**

Total # samples possible =  $(4)(4) = 4^2 = N^n = 16$

Probability of each sample =  $\frac{1}{N^n} = \frac{1}{16}$

Here are the 16 possible (ordered) samples:

(spade, spade) (club, spade) (heart, spade) (diamond, spade)	(spade, club) (club, club) (heart, club) (diamond, club)	(spade, heart) (club, heart) (heart, heart) (diamond, heart)	(spade, diamond) (club, diamond) (heart, diamond) (diamond, diamond)



### What if the Order of the Sample Doesn't Matter?

*Tip! – We will see this again when we learn the **Binomial Distribution**...*

**Ordered Samples** - In the previous examples, the sample obtained was defined by BOTH its membership (eg – A B C) and the ORDER of its members (eg – A first, B second, C third).

**Unordered Samples** – In an Unordered sample, the sample is defined ONLY by its membership.

**Example of an Unordered Sample = “A and B and C”.** There are 6 “qualifying” ordered sequences:

{ A B C }	is the same as
{ A C B }	is the same as
{ B A C }	is the same as
{ B C A }	is the same as
{ C A B }	is the same as
{ C B A }	

What if we just want to know the number of “qualifying” ordered sequences of “A and B and C” that satisfy the event of “A and B and C”? **We don't really want to have to list them all out every time; how tedious.**

The answer is obtained by answering the question: *How many **rearrangements of the orderings** are there of “A”, “B” and “C”?* This is the same as asking: *How many different **permutations** are there of “A”, “B” and “C”?*

**Solution:**

- **First position:**  
From among my selected, what letter should I put into the 1<sup>st</sup> position?  
How many choices are possible?  
Answer = 3
- **Second position, given 1<sup>st</sup> is filled:**  
Having “filled” position 1, from the remainder of my selected, what letter should I put into the 2<sup>nd</sup> position? How many choices are possible?  
Answer =  $(3 - 1) = 2$ .
- **Third position, given 1<sup>st</sup> and 2<sup>nd</sup> filled:** And so on ... How many choices are possible?  
Answer =  $(3-2) = 1$ .
- **Thus, the total number of **ordered rearrangements (permutations)** of 3 items**  
**=  $(3) (3-1) (3-2) = 6$**

How many different ordered rearrangements (permutations) are there of “n” items, such as “A”, “B”, …, “n”?

- **First position:** # choices =  $n$
- **Second position, given 1<sup>st</sup> is filled:** # choices =  $(n - 1)$ .
- **Third position, given 1<sup>st</sup> and 2<sup>nd</sup> are filled:** # choices =  $(n - 2)$ ,

Etc for 4<sup>th</sup> position, 5<sup>th</sup> position and so on to the nth position ....

- **Answer:** The number of permutations of n items is  $= (n)(n-1)(n-2) \cdots (2)(1)$

The number of *ordered* rearrangements (**permutations**) of n items is

$$n! = (n) (n-1) (n-2) \dots (2) (1)$$

**Note** –  $n!$  (called the factorial) is just an abbreviation, so that we don't have to write out long hand expressions like  $(n)(n-1)(n-2)\dots(2)(1)$ . More on this in Unit 6.

**How many equally likely samples are there?**  
**Simple Random Sampling WITHOUT Replacement**  
**Order Does NOT Matter**

**Example – You're playing cards on a Friday night. The dealer has just the 4 queens. He/she gives you 2 of them. What are the chances that your hand contains the queen of clubs and the queen of hearts, regardless of which you got first and which you got second?**

**Population**

Four Queens (N=4)

**Sampling Plan**

- Draw one card at random
- Note its suit
- *Set the selected card aside (Do NOT return it to the pile)*
- Draw a 2<sup>nd</sup> card at random from the 3 that are remaining.
- Note its suit

**Population size, N=4**

**Sample size, n=2**

Total # (ordered sequence) samples possible =  $(4)(4-1) = (4)(3) = 12$

Probability of each (ordered sequence) sample =  $\frac{1}{N(N-1)} = \frac{1}{(4)(3)} = \frac{1}{12}$

**Here is the sample space comprised of all 12 ordered sequence sample points.**

**Under simple random sampling each occurs with equal probability = 1/(12):**

(spade, club) (club, spade)	(heart, club) (club, heart)	(diamond, heart) (heart, diamond)
(heart, spade) (spade, heart)	(diamond, club) (club, diamond)	(spade, diamond) (diamond, spade)

**But, if order DOES NOT matter, then the events of interest are 6 equally likely “hands”**

Hand (Event)	Spade and Club	Heart and Club	Diamond and Heart	Spade and Heart	Diamond and Club	Spade and Diamond
<b>Qualifying Ordered sequences</b>	[ spade, club] [ club, spade]	[ heart, club] [ club, heart]	[ diamond heart] [heart, diamond]	[ spade, heart] [ heart, spade]	[diamond, club] [ club, diamond]	[ spade, diamond] [ diamond, spade]

Nature ——— Population/  
Sample ——— Observation/  
Data ——— Relationships/  
Modeling ——— Analysis/  
Synthesis

**Putting it all together: How to Calculate the Probability of a Sample Outcome, in which order does not matter, under Simple Random Sampling**

The solution involves two calculations of # orderings (permutations)

**(1) Calculate the total number of “ordered sequence” samples**

$$\begin{aligned} \binom{\text{Total \#}}{\text{ordered sequence samples}} &= \binom{\# \text{ ways to}}{\text{draw 1st}} \binom{\# \text{ ways to}}{\text{draw 2nd}} \\ &= (N)(N-1) \\ &= (4)(3) \\ &= 12 \end{aligned}$$

**(2) Next, calculate the *number of ordered rearrangements (permutations) of the sample obtained.***

$$\begin{aligned} \binom{\# \text{ orderings of}}{\text{given sample}} &= \binom{\# \text{ choices for}}{\text{position 1}} \binom{\# \text{ choices for}}{\text{position 2}} \\ &= (n)(n-1) \\ &= (2)(1) \\ &= 2 \end{aligned}$$

**Pr[each equally likely UNordered sample]**

$$= \frac{\# \text{ ordered rearrangements (permutations) of the sample obtained}}{\# \text{ of ordered samples that could have been obtained}}$$

$$= \frac{(n)(n-1)(n-2) \dots (2)(1)}{(N)(N-1)(N-2) \dots (N-n+1)}$$

$$\binom{\text{Probability of}}{\text{1 club and 1 heart}} = \frac{\# \text{ permutations of } n=2 \text{ cards}}{\# \text{ of ordered samples of } n=2 \text{ from } N=4} = \frac{(n)(n-1)}{(N)(N-1)} = \frac{(2)(1)}{(4)(3)} = \frac{2}{12}$$

## Simple Random Sampling **WITHOUT** Replacement Summary

### IF Order **DOES** Matter

- ◆ Total # ordered sequence samples =  $N(N-1)(N-2) \cdots (N-n+1)$
- ◆ Under simple random sampling,

$$\begin{aligned} \text{Probability [ each ordered sequence sample ]} \\ = \frac{1}{(N)(N-1)\dots(N-n+1)} \end{aligned}$$

### IF Order Does **NOT** Matter

- ◆ Total # ordered samples =  $N(N-1)(N-2) \cdots (N-n+1)$
- ◆ # re-arrangements of the sample obtained =  $(n)(n-1)(n-2)\dots(2)(1)$
- ◆ Under simple random sampling,

$$\begin{aligned} \text{Probability [ each sample, regardless of order ]} \\ = \frac{(n)(n-1)(n-2)\dots(2)(1)}{(N)(N-1)\dots(N-n+1)} \end{aligned}$$

- ◆ The total # of Unordered samples is the reciprocal of this.

$$\text{Total \# UNordered samples} = \frac{(N)(N-1)(N-2)\dots(N-n+1)}{(n)(n-1)(n-2)\dots(2)(1)}$$

**b. How to Select a Simple Random Sample WITHOUT Replacement**  
*(Using a random number table)*

Note – You won't ever have to do this in actuality; this is just to illustrate the concept!

<p><b><u>Step 1:</u></b></p> <p>List the subjects in the sampled population.</p> <ul style="list-style-type: none"> <li>♣ This is the <u>sampling frame</u>.</li> </ul>	<p><b><u>Example – Obtain a simple random sample of n=30 from a population of size N=500-</u></b></p> <ul style="list-style-type: none"> <li>♣ Make a list of all N=500</li> </ul>
<p><b><u>Step 2:</u></b></p> <p>Number this listing from “1” to “N”</p> <ul style="list-style-type: none"> <li>♣ where N = size of sampled population</li> </ul>	<p><b><u>Example, continued -</u></b></p> <p style="text-align: center;">N = 500 n = 30</p>
<p><b><u>Step 3:</u></b></p> <p>The size of “N” tells you how many digits in a random number to be looking at:</p> <ul style="list-style-type: none"> <li>♣ For N <math>\leq</math> 10 Need only read 1 digit</li> <li>♣ For N between 10 and 99 Read 2 digits</li> <li>♣ For N between 100 and 999 Read 3 digits</li> <li>etc</li> </ul>	<p><b><u>Example, continued –</u></b></p> <p>Since N=500 is between 100 and 999 and is 3 digits long.</p> <ul style="list-style-type: none"> <li>♣ Read 3 digits</li> </ul>

**Step 4:**

Using the random number table, pick a random number as a starting point

79889	75532	28704
48895	11196	34335
89604	41372	10837



**Example, continued** -

The first 3 digits of this number is “111”. So we will include the 111<sup>th</sup> subject in our sample

**Step 5:**

Proceed down your selected column of the random number table, row by row.

With each row, if the required digits are  $\leq N$ , INCLUDE

With each row, if the required digits are  $> N$ , PASS BY

With each row, if the required digits are a repeat of a previous selection, PASS BY

79889	75532	28704
48895	11196	34335
89604	41372	10837



**Example, continued** -

The first 3 digits of the second random number is “413”. So we will include the 413<sup>th</sup> subject in our sample

**Step 6:**

Repeat “Step 5” a total of  $n=30$  times, which is your desired sample size.

## Remarks on Simple Random Sampling

### Advantages:

- Selection is entirely left to chance.
- Selection bias is still possible, but chances are small.
- No chance for discretion on the part of the investigator or on the part of the interviewers.
- We can compute the probability of observing any one sample. This gives a basis for statistical inference to the population, our ultimate goal.

### Disadvantages:

- We still don't know if a particular sample is representative
- Depending upon the nature of the population being studied, it may be difficult or time-consuming to select a simple random sample.
- An individual sample might have a disproportionate # of skewed values.