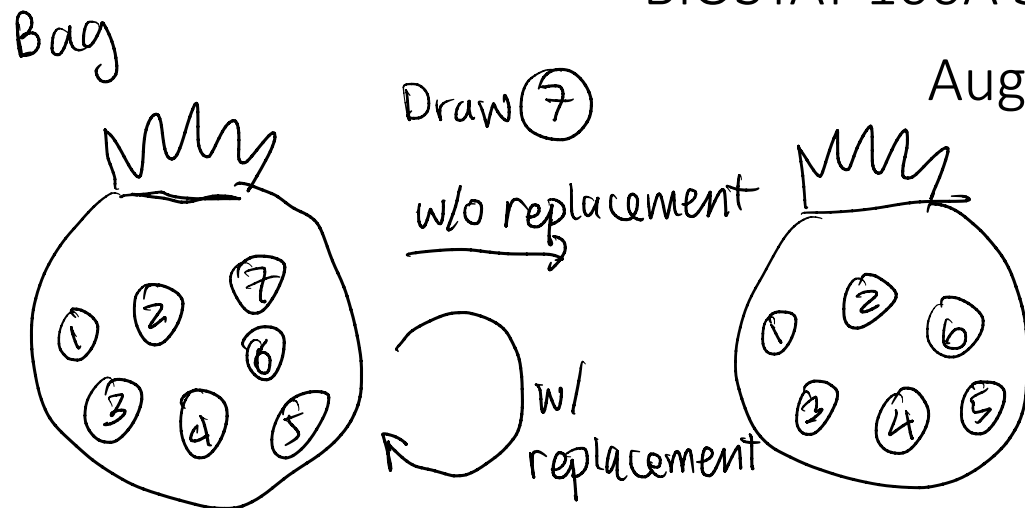


# Week 2 Review

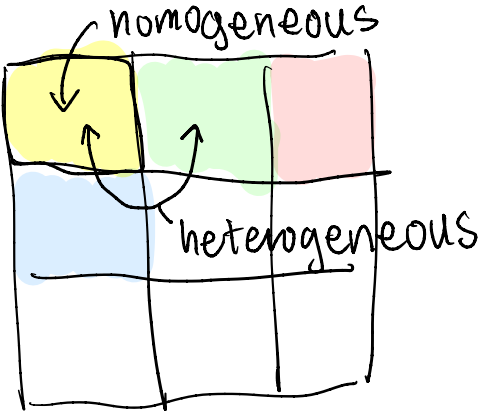
Cindy J. Pang

BIOSTAT 100A Summer Session C 2024

August 16, 2024



# Probability (Random) Sampling

Type of Sampling	When/Why do this type of Sampling ↓ randomization	Selection Mechanism, or How to conduct this type of Sampling
<p>① <u>Simple Random Sampling (SRS)</u> every sample has an equal chance of being selected</p> <p>(1) <u>w/o replacement</u> - we don't put the draw back</p> <p>(2) <u>w/ replacement</u> - put the element back.</p>	<p>convenience</p> <p><u>Assumption</u> - assume respondents in the same are <u>homogeneous</u>. This is <u>problematic</u> b/c most pop'l's are <u>not</u> homogeneous</p>	<p>(1) Population Listing / Pop'n Frame</p> <p>(2) Assign unique ID to each person in the frame</p> <p>(3) Draw randomly.</p>
<p>② <u>Stratified random Sampling</u></p> 	<p><u>When</u>: there are homogeneous subpopulations (strata) / covariates within the pop'n you are interested in</p> <p>⇒ <u>homogeneous within; heterogeneous across.</u></p> <p><math>n_1 = \left(\frac{N_1}{N}\right)n = \left(\frac{30}{100}\right)50 = 15</math> <math>N = \text{pop'n size}</math></p> <p><math>n_2 = \left(\frac{N_2}{N}\right)n = \left(\frac{60}{100}\right)50 = 30</math> <math>n = \text{sample size}</math></p>	<p>÷ the pop'n into strata and take a SRS of each strata.</p> <p><u>Ex: Proportional Allocation</u> - how to select the # of people to include in a strata</p> <p>Suppose you have pop'n w/ 3 strata = <math>\{G_1, G_2, G_3\}</math> where <math>N = 100</math> and <math>N_1 + N_2 + N_3 = N</math>, but we can only sample 50 people</p>

$$n_3 = \left( \frac{N_3}{N} \right) n = \left( \frac{10}{100} \right) 50 = 5$$

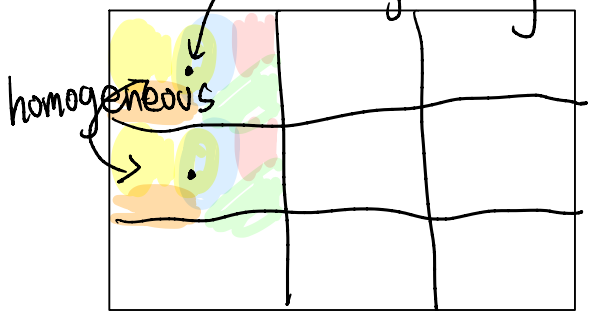
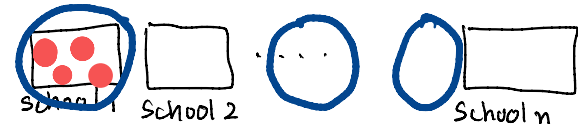
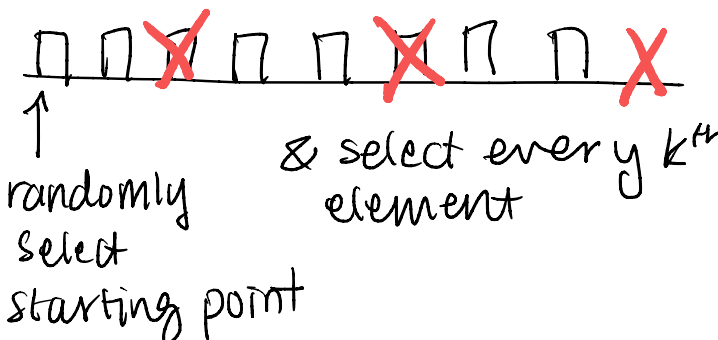
$$N = 100$$

$G_1$	$G_2$	$G_3$
$N_1=30$	$N_2=60$	$N_3=10$
$n_1=15$	$n_2=30$	$5$

$$n_i = \left( \frac{N_i}{N} \right) \cdot n \quad n = 50$$

$i=1, \dots, \# \text{ of strata.}$

# Probability (Random) Sampling

Type of Sampling	When/Why do this type of Sampling	Selection Mechanism, or How to conduct this type of Sampling
<p>③ Cluster Random Sampling</p> <p>"the opposite" of stratified heterogeneity</p> 	<p>When: 1) convenient, ↓ costs 2) Data has natural groups or "clusters"</p> <p>Ex: Test Performance across CA</p>  <p>1-Stage: Take SRS of schools and sample all students in the school</p> <p>2-Stage: After selecting schools, take SRS of teachers within selected schools and sample all their students</p>	<p>1-Stage Use SRS to select clusters and then sample all elements within the cluster</p> <p>2-Stage 1-Stage &amp; SRS to select subgroups within the cluster; take elements in subgroup.</p>
<p>④ Systematic Random Sampling (Line Sampling)</p>	<ul style="list-style-type: none"> <li>Population is dynamic</li> <li>Pop'n frame is in line format ex: list of addresses, phone #s, etc.</li> </ul>	 <p>&amp; select every <math>k^{\text{th}}</math> element</p>

# Data Display

- We can estimate the **Frequency Distribution** with:

- (1) Tables and Graphs
    - Frequency Table
    - Histogram ("Bar Graph")
      - Information from a histogram
    - Cumulative Frequency Polygon
      - percentiles
    - Boxplot
      - rank statistics
      - parts of a boxplot
      - skewness, right vs left skew
  - (2) "Theoretical" Description
    - Normal (Gaussian) Distribution
      - Why is this distribution useful?
    - Log-Normal Distribution
      - Why is this distribution useful?
    - Exponential Distribution
      - When is this distribution useful?
  - (3) Numerical (next lecture)
- Sensitivity vs Specificity
    - Trade-off** between Sensitivity and Specificity → What happens when you move the line?
  - Outliers
    - how to identify outliers
    - what do you do about outliers?

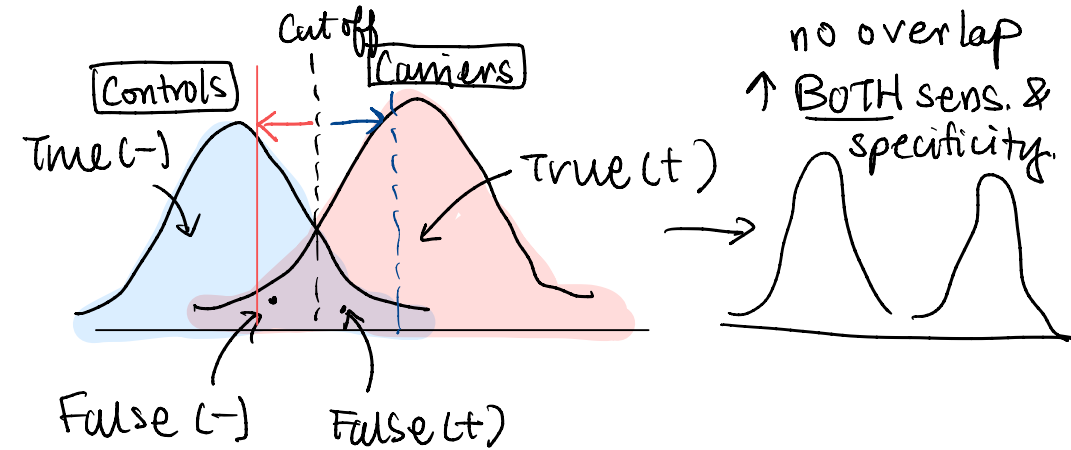
$$\text{Sensitivity} = \frac{TP}{FN + TP} \times 100 \rightarrow \text{"\% of people who correctly diagnosed if they have the disease"}$$

$$= \frac{P(\text{Disorder} \cap \text{Test} +)}{P(\text{Have Dis.})}$$

$$\text{Specificity} = \frac{TN}{FP + TN} \rightarrow \text{"\% of people correctly diagnosed if they don't have the disease"}$$

$$= \frac{P(- \cap D)}{P(-)}$$

	(D) Disorder	( $\bar{D}$ ) No Dis.
(+) Test	True (+) (TP)	False (+) (FP)
(-) Test	False (-) (FN)	True (-) (TN)

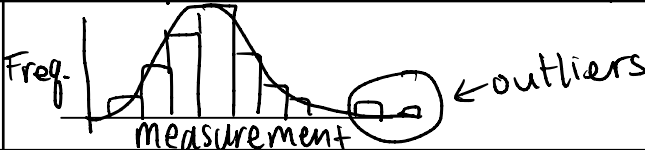
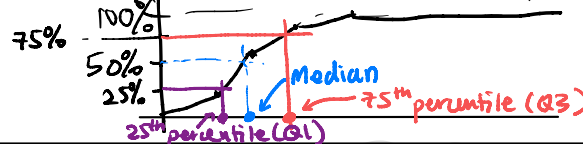
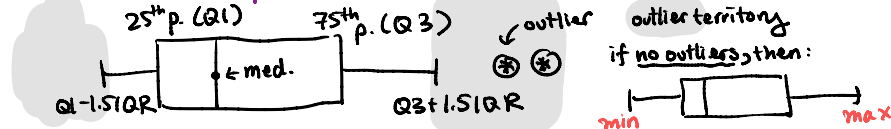
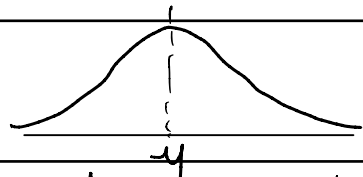
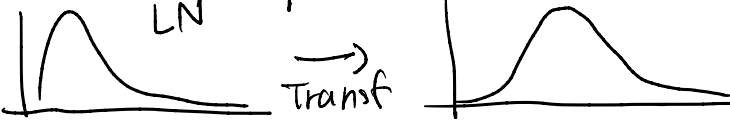



← FP ↑, TN ↓  
Specificity ↓, Sensitivity ↑  
→ FN ↑, TP ↓, ↓ FP  
specificity ↑, sensitivity ↓

↓  
Inverse Relationship




# Descriptive Statistics

Tables and Graphs						
	What it looks like				What does it tell us / Utility?	
	<div>countsPercentageCumulative</div> <div>Last Row = total # of obs - Last Row = 100%</div>					
Frequency Tables	Interval	Abs. Freq.	Rel. Freq (%)	Abs. Freq.	Rel. Freq.	Gives you the numbers
	1-10 11-20 ⋮	5 6 ⋮	1-2% 6-7% ⋮	5 11 ⋮	1-2% 7-9% ⋮	
Histogram						<ul style="list-style-type: none"><li>Shape of our frequency distr.</li><li>outliers → observations that appear "extreme"</li></ul>
Cumulative Frequency Polygon						Estimating <u>Percentiles</u>
Boxplot						<ul style="list-style-type: none"><li>shape</li><li>identify outliers (outside whiskers)</li></ul>
Theoretical Descriptions						
Normal Distribution	 <div><math>N(\mu, \sigma^2)</math> mean↑variance</div>					<ul style="list-style-type: none"><li>symmetric data</li><li>mean = median</li></ul>
Log-Normal Distribution						<ul style="list-style-type: none"><li>skewed data (right skewed)</li></ul>
Exponential Distribution						survival data

# Numerical Descriptions of Data

- Measures of Location
  - (1) Arithmetic Mean (average)
  - (2) Median – how to find the median when it is even vs odd
  - (3) Geometric Mean
  - (4) Mode
  - (5) Midrange

## Measure of Location

	Formula/ How to calculate it	When to use it	Statistic → Sample	Parameter → Population
Arithmetic Mean	$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{x_1 + \dots + x_n}{n}$	symmetric distr.	$\bar{x}$ = sample mean	$\mu$ = pop'n mean
Median	order all values → $\frac{1}{2} < 50\%$ percentile $\frac{1}{2} > 50\%$ percentile	skewed.		
Geometric Mean	1.) $\log(x_i)$ , $i=1, \dots, n$ 2.) mean of the logs $\bar{x}_{\log} = \frac{1}{n} \sum_{i=1}^n \log(x_i)$	3.) $10^{\bar{x}_{\log}}$ exponential distr.		
Mode	most freq. value	skewed data, bimodal   ← no mode		
Midrange	$\frac{\text{Max} + \text{Min}}{2}$	Quick & Dirty Statistic		