SAMPLING DESIGN

- Simple Random Sample
- Systematic Sample
- Stratified Random Sample
- Cluster Sample
- Bad Sampling Methods
- Errors in Data Collection



MOTIVATION: SAMPLING PLANS

- Scenario: Allegheny County has 974 bridges that are at least 20 feet long. Our goal is to learn about the lengths of the bridges.
 - Dataset: Contains 974 rows, each with the municipality, rating, and length
- Question: Why is this dataset a population?
- Answer: ______ in Allegheny County is included
 - Allows us to calculate ______ for length
- Scenario: Suppose we only want to work with a subset of the data.
- Question: How can we be confident that the sample we obtain will give us accurate information about the population?
- Answer: Use a ______

PROBABILITY SAMPLING PLAN

- **Probability sampling plan:** a method of selecting a sample from a population that utilizes some form of random selection
 - A good plan will return a sample that is **representative** the sample resembles the population across all characteristics
 - **Example:** If a population is half male and half female, the sample will have about half males and half females as well
 - **Example:** If the average age in a population is 40, then the average age in the sample is also around 40.
 - Several different types
 - Simple random sample
 - Systematic sample
 - Stratified random sample
 - Cluster sample

SIMPLE RANDOM SAMPLE

- Simple random sample: a probability sampling plan where every member of the population has an equal chance of being selected for inclusion in the sample
 - Equivalent to putting every observation's name in a hat and randomly pulling out names without replacing them to obtain the sample

1	2	3	4	5	6	7	8	9	10
11	12	13	14	15	16	17	18	19	20
21	22	23	24	25	26	27	28	29	30
31	32	33	34	35	36	37	38	39	40
41	42	43	44	45	46	47	48	49	50

Example: Sampling Methods in Action

- Scenario: Data contains every bridge in Allegheny County.
- Task: Use a simple random sample to estimate the average bridge length and compare it against the average length of all bridges
- Strategy: Use R to select a _____
 - Read in the dataset
 - Use the **sample** function to select observations and create a subset
 - Choose 50 observations
 - Average the bridge lengths in the subset to obtain the statistic
- Parameter: $\mu =$ _____
- Statistic: $\bar{x} =$ _____

OTHER PROBABILITY SAMPLING METHODS

- Systematic sample: members are sampled according to some predetermined rule by skipping a certain number of subjects and then sampling every n^{th} observation
- Stratified random sample: two step process where:
 - Population is first divided into groups of similar individuals (called strata)
 - Simple random sample is chosen from each stratum and combined to form the full sample
- Cluster sample: three step process where:
 - Population is first divided into groups (called clusters)
 - Simple random sample of the clusters is taken
 - Every member in the selected clusters becomes part of the sample

4

EXAMPLE: SYSTEMATIC SAMPLE
Scenario: Data contains every bridge in Allegheny County.
 Task: Use a systematic sample to estimate the average bridge length.
Strategy:
 Select every observation for inclusion in the sample Includes observation numbers,,,, Total of observations Average sampled observations
• Parameter: $\mu = 284.24$
• Statistic: $\bar{x} =$
FXAMPLE COMPARING STRATIFIED AND CLUSTER SAMPLING
Construction of Supreme Court for
Democrats, Republicans, and Independents.
• Ouestion: What sampling method should be used?
• Answer:
From a list of registered voters, randomly sample a subset of people from
 Eind who support the Supreme Court for each affiliation and
Expect there to be
EXAMPLE COMPARING STRATIFIED AND CLUSTER SAMPLING
• Scenario: Want to estimate the average water bill of households in
a neighborhood with 20 streets by going door to door and asking residents what their bill was last month.
 Question: What sampling method should be used?
• Answer:
Streets in the same neighborhood should be so a random sample is not necessary
Randomly choose a (say 3) and visit on these streets
Better than a random sample where a from streets might be selected



• **Convenience sample:** a sample obtained from members of the population who were the easiest to access

EXAMPLE: VOLUNTARY SAMPLE
 Scenario: After a political debate, CNN wants to gauge if its viewers believed the Democratic or Republican candidate won the debate. They put a link to a poll on CNN.com across the scroll at the bottom of screen and invite the viewing audience to respond. Question: Why is this a voluntary sample? Answer: Respondents and no sampling is done
• Ouestion: Why is this likely to be biased?
 Answer: Only viewers who one way or another will take time to vote CNN has a more so sample is likely to vote for the as the winner
EXAMPLE: CONVENIENCE SAMPLE
 Scenario: A new store is considering moving into a shopping mall. A representative for the store polls the first 100 shoppers she encountered on how likely they would be to shop at this store. Question: Why is this a convenience sample? Answer: Responses taken from people who were
• Ouestion: Why is this sample likely to be biased?
 Answer: Likely sampled groups who came to the mall who have opinions Only was used Sample taken on (many people) is very different from sample taken on (many people and parents)
ERRORS IN DATA COLLECTION

- Sampling error: refers to the difference between the statistic and parameter due to the observations selected for the sample
 - Occur in every sample because the statistic (almost) never equals the parameter
 - Despite the name, they are not mistakes and occur naturally
- Nonsampling error: occur due to mistakes that occur in the sampling process, either due to human error or other problems
 - Data acquisition: recording data or taking measurements incorrectly
 - Nonresponse bias: observation selected for the sample does not respond
 - Selection bias: some observations in the population cannot be chosen for inclusion in the sample

EXA	MPLE: SAMPLING ERROR	16
• Sce 974 rand	nario: The parameter for the length (i.e. average length of all bridges) was 284.24, while the statistic from the simple dom sample (i.e. average length of 50 bridges) was 296.72.	
• Que	stion: What is the sampling error?	
• Ans	wer: (Difference between and	_)
• Que • Ans • S	stion: Should we be concerned about this difference? wer: amples are – only "missed" the parameter by	
Exa	MPLE: Nonsampling Errors	17
EXA • Tas	MPLE: NONSAMPLING ERRORS	17
• Tasl	MPLE: NONSAMPLING ERRORS C Match each of the following nonsampling errors with the esponding type.	17
• Tasl corr	MPLE: NONSAMPLING ERRORS K: Match each of the following nonsampling errors with the esponding type. A simple random sample is chosen from a list of residents who have a driver's license.	17
EXA • Tasl corr I. II.	MPLE: NONSAMPLING ERRORS C Match each of the following nonsampling errors with the esponding type. A simple random sample is chosen from a list of residents who have a driver's license. A scale in a laboratory is not calibrated correctly and adds 3 grams to each measurement.	17
• Tasl corr I. II.	MPLE: NONSAMPLING ERRORS C Match each of the following nonsampling errors with the esponding type. A simple random sample is chosen from a list of residents who have a driver's license. A scale in a laboratory is not calibrated correctly and adds 3 grams to each measurement. A subject ignores a phone call from a pollster during election season.	17
EXAI • Tasl corr I. III.	MPLE: NONSAMPLING ERRORS & Match each of the following nonsampling errors with the esponding type. A simple random sample is chosen from a list of residents who have a driver's license. A scale in a laboratory is not calibrated correctly and adds 3 grams to each measurement. A subject ignores a phone call from a pollster during election season.	17
EXAI • Tasl corr I. II. III.	MPLE: NONSAMPLING ERRORS A Simple random sample is chosen from a list of residents who have a driver's license. A scale in a laboratory is not calibrated correctly and adds 3 grams to each measurement. A subject ignores a phone call from a pollster during election season. wer:	17
EXAI • Tasl corr I. II. III. • Ans I.	MPLE: NONSAMPLING ERRORS	17
EXAI • Tasl corr I. II. III. III. III.	MPLE: NONSAMPLING ERRORS	17