# Unit 3 Kickoff Packet

*Collecting Data: Sampling & Study Design*

Topics 3.1, 3.2, and 3.3 — Video Follow-Along Assignment

---

### Learning Objectives for This Unit

- **VAR-1.E:** Identify questions to be answered about data collection methods.
- **DAT-2.A:** Identify the type of a study (observational vs. experiment).
- **DAT-2.B:** Identify appropriate generalizations and conclusions based on study type.
- **DAT-2.C:** Identify a sampling method, given a description of a study.
- **DAT-2.D:** Explain why a particular sampling method is or is not appropriate.

## PART 1: Unit 3 Cheat Sheet — Key Terms & Definitions

| Term | Definition |
|---|---|
| **Population** | The *entire group* of individuals we want information about. Think: population = *all.* |
| **Sample** | A *subset* of the population from which we actually collect data. |
| **Census** | Collecting data from *every* individual in the population (often impractical). |
| **Parameter** | A number that describes the *population* (often unknown). |
| **Statistic** | A number that describes the *sample* (used to estimate parameters). |

| | Sampling Methods |
|---|---|
| **Simple Random Sample (SRS)** | Every group of $n$ individuals has an *equal chance* of being selected. Use a random number generator. |
| **Stratified Sample** | Divide population into *homogeneous* groups (strata), then take an SRS *within each* group. |
| **Cluster Sample** | Divide population into *heterogeneous* groups (clusters), then randomly select *entire groups* and sample *all* individuals in chosen groups. |
| **Systematic Sample** | Randomly choose a starting point, then sample at *fixed intervals* (e.g., every 20th person). |

| | Study Types & Conclusions |
|---|---|
| **Observational Study** | No treatments are imposed; we just *observe*. **Cannot determine cause-and-effect.** |
| **Experiment** | Treatments *are imposed* on subjects. **Can determine cause-and-effect** (if well-designed). |
| **Bias** | Systematic error that causes estimates to be *inaccurate* (not centered at truth). |
| **Variability** | How *spread out* estimates are; low variability = high *precision*. |

## Common Confusion: Cluster vs. Stratified

| CLUSTER | STRATIFIED |
|---|---|
| • Ideal groups are **heterogeneous** (mixed) | • Ideal groups are **homogeneous** (similar) |
| • Select random **groups** | • Sample **within each** group |
| • Sample **ALL** in selected groups | • Sample **SOME** from every group |
| • Easier but can have high variability | • More precise but harder to implement |

# PART 2: Guided Video Notes

## Video 1: Topic 3.1 — Do the Data Tell the Truth?  [∼5 min]

> **Key Concept:** Methods for data collection that do not rely on chance can result in untrustworthy conclusions. The way data is collected *completely informs* how we analyze it.

### The Abraham Wald Airplane Problem  [0:30–3:30]

1. Abraham Wald worked for the _____, helping the _____ forces during World War II through mathematical analysis.

2. The British military had a problem: bombers were being shot down over Nazi Germany. They wanted to add heavy _____ to reinforce planes, but could only put it in _____ place.

3. Soldiers made bullet hole charts for planes that _____ from bombing missions and asked Wald where to add armor.

4. **Pause & Think:** Looking at the bullet hole patterns, where would most people initially suggest putting armor? _____

5. **Wald's Key Insight:** The critical realization was about the _____ in which the data was collected.

6. The planes in the sample were only the ones that _____. The bullet holes we see mark places where planes can take hits and still _____!

7. Therefore, armor should go where we *don't* see bullet holes (like the _____), because planes hit there never made it back.

---

### Population vs. Sample in the Wald Problem

| | |
|---|---|
| **Population** (what we want to know about) | All planes that were hit during bombing missions |
| **Sample** (what we actually measured) | Only planes that made it back |
| **Who was missing?** | Planes that were shot down |
| **Type of Bias** | **Survivorship Bias** |

---

### Key Takeaways from Video 1  [4:50–5:12]

8. A proper analysis of data must take into account _____.

9. Sometimes, our samples may not be _____ of the whole _____.

10. The statistician's motto: "Be _____, Be _____, Be _____. Avoid BS (_____)."

# Video 2: Topic 3.2 — Planning a Study              [∼8 min]

> **Key Concept:** The type of study determines what conclusions we can draw. Observational studies *cannot* establish causation—only experiments can.

## Generalization Rules                                   [3:17–3:57]

1. It is only appropriate to generalize about a population based on samples that are:

   - _____ or otherwise representative, AND
   - Selected from *that specific* _____

2. Example: If we study lima beans in humid climates, we _____ generalize to black beans, because they are a different population.

## The Racial Income Gap: An Observational Study            [3:57–6:10]

3. The 2018 Current Population Survey found:

   - White households median: $70,642      Black households median: $58,665
   - Difference: ≈ $_____/year → $_____ over a 45-year career

4. Possible causes include: inequity of schools, familial connections, or direct discrimination. These are called _____ variables.

5. This survey is an _____ because no treatments were imposed. Therefore, we _____ infer cause and effect.

6. A **retrospective** study examines _____ data; a **prospective** study follows individuals into the _____.

## The Résumé Experiment                                  [6:10–7:47]

7. Researchers sent identical résumés to employers, _____ assigning either a commonly white or commonly Black name.

8. They measured the _____ rates for both groups.

9. This is an _____ because treatments (name type) are _____ upon subjects.

10. If well-designed, an experiment _____ determine a causal relationship.

## Fill In: Observational Study vs. Experiment

|  | Observational Study | Experiment |
|---|---|---|
| **Treatments imposed?** | _____ | _____ |
| **Can show causation?** | _____ | _____ |
| **Example from video** | _____ | _____ |

# Video 3: Topic 3.3 (Part 1) — Random Sampling Methods    [∼9 min]

> **Context:** San Antonio is one of the most economically **segregated** cities in the U.S.—incomes are similar *within* neighborhoods but vary greatly *between* neighborhoods. We want to estimate the median household income using different sampling methods.

## Why Sample?        [0:00–2:24]

1. A **census** collects data from _____ individuals, but censuses are _____ to conduct.

2. Due to the right-skewed nature of income data, the _____ is preferred over the mean.

3. If done well, a random sample should be _____ of the general population.

## Simple Random Sample (SRS)        [2:24–5:00]

4. An SRS is a sample where every _____ of size $n$ has an _____ chance of being chosen.

5. **How to conduct:** Number all individuals 1 to $N$, then use a _____ to select $n$ numbers _____ (don't repeat).

6. The SRS sample median was **$50,500**.

## Cluster Random Sample        [5:00–6:13]

7. Divide the population into _____ that are near one another.

8. Take an SRS of entire _____, then sample _____ individuals in selected clusters.

9. In San Antonio, 2 clusters were randomly selected out of 100 regions. The cluster sample median was **$110,350**.

10. Advantage: Easier to collect (only need to visit _____ areas).

## Stratified Random Sample        [6:13–7:20]

11. Divide the population into _____ based on a similar characteristic.

12. Take an SRS _____ each stratum, then combine all selected individuals.

13. In San Antonio, 100 homes were randomly selected from *each* of the 100 regions. The stratified sample median was **$51,025**.

**Systematic Random Sample**                                              [8:09–8:48]

14. Randomly choose a _____, then sample at a fixed _____ interval.

15. Example: In a lunch line, randomly select a number 1–20, then sample every _____ person.

16. Advantage: Very _____, especially when individuals are "lined up."

# Video 4: Topic 3.3 (Part 2) — Evaluating Sampling Methods   [∼8 min]

> **Key Question:** The three methods gave different estimates ($50,500, $110,350, $51,025). Which should we trust? We evaluate using **bias** (accuracy) and **variability** (precision).

## Understanding Bias & Variability          [1:42–2:52]

1. **Bias** measures _____: Are estimates centered at the _____?

2. **Variability** measures _____: How spread out are the estimates?

3. The "gold standard" is _____ bias and _____ variability.

4. The *true* median household income in San Antonio (from the census) is $_____.

---

**Label the Dartboards**

*Match each description: Biased + High Var, Biased + Low Var, Unbiased + High Var, Unbiased + Low Var*

| Target A | Target B | Target C | Target D |
|---|---|---|---|
| (off-center, spread) | (off-center, clustered) | (centered, spread) | (centered, clustered) |
| _____ | _____ | _____ | _____ |
| _____ | _____ | _____ | _____ |

---

## Simulation Results for Each Method          [2:52–7:42]

5. **Non-Random Sample** (students picking "representative" areas):

   - Result: _____—students systematically overestimated.
   - Lesson: Non-random samples can lead to _____.

6. **Simple Random Sample (SRS):**

   - Bias: _____ (centered at true median)
   - Variability: _____
   - Advantage: Unbiased, easy to _____
   - Disadvantage: Can be difficult to _____; may not be as precise

7. **Cluster Random Sample:**

   - Bias: _____ (uses random selection)
   - Variability: _____ (estimates spread widely)

- Why? Income is homogeneous _____ regions but varies _____ regions.
- Works best when clusters are _____ and similar to one another.

8. **Stratified Random Sample:**

   - Bias: _____
   - Variability: _____ (very precise!)
   - Why? Each sample has a similar _____ of incomes.
   - Works best when strata are _____.
   - Disadvantage: Can be _____ to implement.

---

### Fill In: Sampling Methods Comparison

| Method | Bias? | Variability | Best When... |
|---|---|---|---|
| Simple Random Sample | _____ | _____ | _____ |
| Cluster Sample | _____ | _____ | Clusters are _____ |
| Stratified Sample | _____ | _____ | Strata are _____ |

---

## Key Takeaways from Video 4        [7:42–7:58]

9. Random sampling tends to provide _____ estimates.

10. Cluster sampling is effective when clusters are _____ and _____ to each other.

11. Stratified sampling is most effective when strata are _____.

# PART 3: Check for Understanding

> **Directions:** Answer the following questions to check your understanding of Topics 3.1–3.3. Show your reasoning!

## Section A: Observational Studies vs. Experiments

**A1.** An observational study found that sleep is associated with job performance ($r = 0.86$). A reader concluded that more sleep *causes* better performance. Why is this conclusion **not** correct?

    (A) The correlation value should equal 1 for such a conclusion.

    (B) The correlation value should be negative for such a conclusion.

    (C) The sample was not representative of the population.

    (D) Causation cannot be determined from an observational study.

    (E) The correlation value implies less than 75% of variability is explained.

    **Answer:** _____      **Explain why:** _____

**A2.** A researcher studies the effects of an herbal supplement on colds. From 50 people with colds, 25 are assigned to take the supplement and 25 are asked to drink water. The duration of each cold is recorded. What are the **experimental units**?

    (A) All people with a cold

    (B) The sample of 50 people who had a cold

    (C) The 25 people who were given the supplement

    (D) The 25 people who drank water

    (E) The recorded number of days that the cold lasted

    **Answer:** _____

**A3.** A researcher selects a simple random sample of 1,200 women who are students at Midwestern colleges. To which population can results be generalized?

    (A) All students in the United States

    (B) All college students in the United States

    (C) All women who are students in the United States

    (D) All students at Midwestern colleges in the United States

    (E) All women who are students at Midwestern colleges in the United States

    **Answer:** _____      **Why?** _____

## Section B: Sampling Methods (Focus on Cluster vs. Stratified)

**B1.** At a clothing store, clothes are displayed on racks. Clothes on each rack have *similar prices*, but prices *vary greatly between racks*. A consumer randomly selects 4 pieces from *each* rack to estimate the typical price. What type of sample is this?

    (A) A census

    (B) A cluster sample

    (C) A simple random sample

    (D) A stratified random sample

    (E) A systematic random sample

    **Answer:** _____     **Explain:** _____

**B2.** A school district has 15 high schools. Seniors' post-graduation plans *vary greatly from one school to the next*. The superintendent selects 5 schools at random and surveys *all* seniors at those schools. What is a disadvantage of this cluster sample?

    (A) Cluster sampling is usually too expensive.

    (B) The sample will be too large to yield accurate results.

    (C) The schools in the sample might not represent all seniors.

    (D) Absent students on survey day could affect results.

    (E) There is no disadvantage.

    **Answer:** _____     **Explain in terms of bias/variability:**

**B3.** Margo wants to estimate the percent of red marbles in a bag. She randomly selects a marble, records its color, *puts it back*, shakes the bag, and repeats. What type of sampling is this?

    (A) Cluster sampling

    (B) Stratified random sampling

    (C) Systematic random sampling

    (D) Random sampling with replacement

    (E) Random sampling without replacement

    **Answer:** _____

**Section C: You Try It! — Choosing the Right Method**

**C1. Scenario:** A company has offices in 5 different cities. Commute times are *similar within each city* but *vary greatly between cities.* A researcher wants to estimate the average commute time for all employees.

    (a) Should the researcher use **cluster** or **stratified** sampling? Circle one.

    (b) Explain your reasoning using the concepts of homogeneity/heterogeneity:

    (c) What would happen if the researcher used the *wrong* method? (Discuss variability.)

**C2. Scenario:** Wildlife biologists want to estimate the deer population in a region with three distinct areas: a forest, a lake, and a town. Deer populations differ by area.

    (a) Why would **stratified** sampling be better than a simple random sample for placing observation stations?

    (b) If researchers marked deer and recounted them later, what problem arises if the same deer are counted multiple times? How would this affect the population estimate?

**Exit Ticket: Unit 3 Kickoff Summary**

In 2–3 sentences, explain the **two most important ideas** from Unit 3 so far. Include at least one idea about sampling methods and one about study design.