

# 1 Lecture 1: Introduction to Sampling and Distributions - Basic Concepts Before CLT

Duration: 50 minutes

**Learning Outcomes:** By the end of this lecture, students will be able to distinguish between a population and a sample, understand the difference between parameters and statistics, grasp the concept of a sampling distribution of the sample mean, and appreciate the importance of repeated sampling. They will also understand that the sample mean is a random variable and different samples yield different means.

---

## I. Introduction (5 minutes)

- **Welcome and Overview:**

- “Good morning/afternoon, everyone! Today, we’re diving into the exciting world of statistics, specifically focusing on how we gather and interpret data. This lecture lays the groundwork for understanding one of the most powerful theorems in statistics: the Central Limit Theorem (CLT), which we’ll cover in our next session.”
- “Before we get to CLT, we need to understand some fundamental concepts. Think of it as building blocks for a strong foundation.”

## II. Population vs. Sample (10 minutes)

- **What is a Population?**

- “Imagine you want to know the average height of all students in our college. Would it be practical to measure every single student? Probably not!”
- **Definition:** “A population is the entire group of individuals or objects that we are interested in studying. It’s the complete set of all possible observations.”
- **Examples:**
  - \* All students in a university
  - \* All cars manufactured by a specific company in a year
  - \* All trees in a particular forest
  - \* All possible outcomes of rolling a fair die infinitely many times
- “Often, populations are too large or even theoretically infinite, making it impossible to collect data from every member.”

- **What is a Sample?**

- “Since we can’t measure everyone in the college, what do we do? We take a smaller, manageable group.”

- **Definition:** “A sample is a subset or a smaller, representative group drawn from the population. We study the sample to make inferences or draw conclusions about the entire population.”
- **Examples:**
  - \* A group of 100 randomly selected students from the university
  - \* A batch of 50 cars inspected from the production line
  - \* Measurements from 20 trees randomly selected from the forest
- **Analogy:** “Think of it like tasting a spoonful of soup to know if the entire pot needs more salt. The spoonful is your sample, the entire pot is your population.”
- **Key Idea:** “The goal is for the sample to be as representative of the population as possible to ensure our conclusions are accurate.”

### III. Parameters vs. Statistics (10 minutes)

- **Why Differentiate?**

- “Once we have our population and sample, we’ll start calculating values. It’s crucial to distinguish between values that describe the population and values that describe the sample.”

- **Parameters:**

- **Definition:** “A parameter is a numerical characteristic that describes an entire population. It’s usually a fixed, but often unknown, value.”
- **Notation & Examples:**
  - \* Population Mean ( $\mu$  - mu): “This is the true average of a characteristic for the entire population. For example, the actual average height of all students in our college.”
  - \* Population Standard Deviation ( $\sigma$  - sigma): “This measures the spread or variability of data for the entire population. It tells us how much individual data points typically deviate from the population mean.”
- “We use Greek letters for parameters because they are typically unknown and we can only estimate them.”

- **Statistics:**

- **Definition:** “A statistic is a numerical characteristic that describes a sample. It’s calculated from the data in the sample and is used to estimate population parameters.”
- **Notation & Examples:**
  - \* Sample Mean ( $\bar{x}$  - x-bar): “This is the average of a characteristic calculated from our sample. If we measure 100 students, their average height is  $\bar{x}$ .”

- \* Sample Standard Deviation ( $s$ ): “This measures the spread or variability of data within our sample. It’s our best estimate of the population standard deviation based on the sample data.”
- “We use Roman letters for statistics because they are calculated from observed data.”

- **Quick Check:**

- “If I say ‘the average age of all humans on Earth is...’, am I talking about a parameter or a statistic?” (Expected answer: Parameter)
- “If I survey 50 people and find their average income is \$X,000, is that a parameter or a statistic?” (Expected answer: Statistic)

#### IV. Sampling Distributions of the Sample Mean (15 minutes)

- **The Big Idea:** “Now, let’s get to a very important concept: the sampling distribution.”

- “Imagine we take one sample from our college students and calculate their average height ( $\bar{x}_1$ ). Then, we take another sample and get  $\bar{x}_2$ . And another,  $\bar{x}_3$ . Will these sample means be the same?” (Expected answer: No)

- “Exactly! Each time we take a different sample, we’ll likely get a slightly different sample mean. This variation is key.”

- **What is a Sampling Distribution?**

- **Definition:** “The sampling distribution of the sample mean is the probability distribution of all possible sample means that could be drawn from a population for a given sample size ( $n$ ). It tells us how the sample means vary from sample to sample.”
- “It’s not about the distribution of the original data, but the distribution of the means of many, many samples.”

- **Importance of Repeated Sampling:**

- “To understand the sampling distribution, we need to think about taking many samples. This is called repeated sampling.”
- “Why is this important? Because a single sample mean might be high, or low, just by chance. By looking at many sample means, we start to see a pattern – a distribution – of where these means tend to fall.”
- “This pattern will eventually reveal something very powerful about how sample means relate to the true population mean, which leads us to the Central Limit Theorem later.”

## V. Activity: Building a Sampling Distribution (10 minutes)

- **Hands-on Example:** “Let’s make this concrete. We’ll use a very small, hypothetical population to see how this works.”

- **Dataset:** “Our population consists of the following 6 numbers: [5, 10, 15, 20, 25, 30]”

- Optional: Quickly calculate the population mean ( $\mu$ ) and standard deviation ( $\sigma$ ) for demonstration, but emphasize we usually don’t know these.

- $\mu = (5 + 10 + 15 + 20 + 25 + 30)/6 = 105/6 = 17.5$

- **Activity 1: Samples of Size n=2**

1. **Instructions:** “Now, let’s draw all possible samples of size 2 from this population. We’ll do this with replacement for simplicity, meaning we can pick the same number twice. Let’s list them out and calculate the sample mean ( $\bar{x}$ ) for each.”

- (5, 5)  $\rightarrow \bar{x} = 5$

- (5, 10)  $\rightarrow \bar{x} = 7.5$

- (5, 15)  $\rightarrow \bar{x} = 10$

- ...

- (30, 30)  $\rightarrow \bar{x} = 30$

- (There will be  $6 \times 6 = 36$  possible samples of size 2 with replacement. Don’t need to list all, just enough to illustrate the point.)

2. **Calculate Sample Means:** “For each sample, calculate its mean.”

3. **Plot Frequency (Conceptual):** “Now, if we were to create a histogram or a frequency table of all these sample means, what would we see? We’d see some means appear more frequently than others.”

- Illustrate with a quick sketch: “You’d see that sample means tend to cluster around the true population mean (17.5). It would look something like a bell-shaped curve, even with this small dataset.”

- **Activity 2: Samples of Size n=3 (Briefly discuss)**

1. **Instructions:** “What if we took samples of size 3? ( $6 \times 6 \times 6 = 216$  possible samples with replacement). It would be even more combinations.”

2. **Discussion Point:** “What do you think would happen to the spread of the sample means if we increased the sample size from 2 to 3? Would they be more spread out or less spread out around the population mean?”

- (Expected answer: Less spread out. Larger samples tend to be more representative, so their means should be closer to the true population mean.)

## VI. Outcome and Conclusion (5 minutes)

- **Key Takeaways from the Activity:**

- “What did we learn from this exercise?”
- “The sample mean ( $\bar{x}$ ) is a random variable.” “Its value changes from sample to sample, purely due to the random nature of sampling. You don’t get the same  $\bar{x}$  every time.”
- “Different samples give different means.” “This is a fundamental concept. We need to account for this variability when making inferences.”
- “Despite the variability, notice how the sample means tended to cluster around the population mean ( $\mu = 17.5$ ). This is a preview of the Central Limit Theorem!”

- **Recap of Lecture:**

- “Today, we defined populations (the whole group) and samples (a subset).”
- “We distinguished between parameters (population characteristics, like  $\mu$  and  $\sigma$ ) and statistics (sample characteristics, like  $\bar{x}$  and  $s$ ).”
- “Most importantly, we introduced the idea of a sampling distribution of the sample mean – the distribution of all possible sample means. And we saw why repeated sampling helps us understand this distribution.”

- **Looking Ahead:**

- “This understanding of sampling distributions, particularly how sample means behave, is crucial for our next lecture where we will formally introduce the Central Limit Theorem. The CLT tells us exactly what shape this sampling distribution takes and how its variability changes with sample size, even if the original population isn’t normally distributed.”

## Practice Questions & Answers: Lecture 1 - Basic Concepts Before CLT

**Instructions:** Please answer the following questions based on the concepts discussed in today’s lecture.

---

### Question 1: Multiple Choice

Which of the following best describes a population?

- A) A small, representative group selected for study
- B) The entire group of individuals or objects that we are interested in studying
- C) A numerical characteristic calculated from a sample
- D) The average of all possible sample means

**Answer:** B) The entire group of individuals or objects that we are interested in studying.

---

### Question 2: Fill in the Blanks

1. A numerical characteristic that describes an entire population is called a \_\_\_\_\_, while a numerical characteristic calculated from a sample is called a \_\_\_\_\_.
2. We use the symbol  $\mu$  to denote the \_\_\_\_\_ mean, and  $\bar{x}$  to denote the \_\_\_\_\_ mean.

**Answer:**

1. Parameter, Statistic
  2. Population, Sample
- 

### Question 3: True or False

1. The sample mean ( $\bar{x}$ ) is a fixed value that does not change from one sample to another.
2. The sampling distribution of the sample mean is the distribution of the individual data points in the population.
3. Repeated sampling is important to understand how sample means vary.

**Answer:**

1. False (The sample mean is a random variable and changes from sample to sample)
  2. False (It's the distribution of all possible sample means)
  3. True
- 

### Question 4: Short Answer

Imagine a large university wants to determine the average amount of time its students spend studying per week.

1. What would be the population in this scenario?
2. If they survey 500 students and calculate their average study time, what is this calculated average called (a parameter or a statistic)? What symbol would we use for it?

**Answer:**

1. The population would be all students enrolled in the university
  2. This calculated average is a statistic. We would use the symbol  $\bar{x}$  (x-bar) for it
- 

**Question 5: Conceptual Understanding**

You conducted the activity in class with the dataset [5, 10, 15, 20, 25, 30] and calculated sample means for samples of size 2.

1. Why did you get different sample means even though you were drawing from the same population?
2. If you were to take many, many more samples of size 2 and plot all their means, what would you expect the general shape of this plot (the sampling distribution) to look like, even with this small population?

**Answer:**

1. You got different sample means because of sampling variability or random chance in the selection process. Each sample is a different subset of the population, and unless the sample perfectly mirrors the population (which is rare), its mean will likely differ from other sample means. This demonstrates that the sample mean ( $\bar{x}$ ) is a random variable.
2. Even with a small population, if you plot many sample means, the sampling distribution would tend to be bell-shaped and centered around the true population mean ( $\mu$ ). This is an early glimpse of the Central Limit Theorem at play, suggesting that sample means tend to cluster around the population's true average.