

Chapter 15

Sampling Distributions

Intended Learning Outcomes

- Distinguish parameters from statistics
- State the Law of Large Numbers
- Define and interpret a sampling distribution
- Calculate the standard error of \bar{X}_n
- State and apply the Central Limit Theorem
- Compute probabilities involving \bar{X}_n

Guiding Questions

Our goal in statistics is to learn about *populations* using *samples*.

But samples vary, so **how can we make reliable inferences?**

How do we know if an observed result reflects
real evidence or just **random chance**?

PART 1

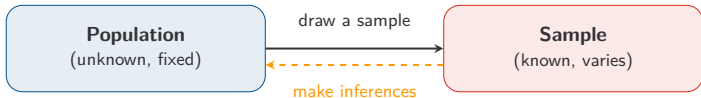
Parameters and Statistics

Distinguishing what we know from what we estimate

Parameters and Statistics

Parameter vs. Statistic

- A **parameter** is a number that describes the entire *population*. Parameters are typically **unknown** and **fixed**.
- A **statistic** is a number computed from *sample* data. Statistics are **known** but **vary** from sample to sample.



Identify Parameters and Statistics

Example 15.1

}

CORRECTION • Parameter vs. Statistic Examples (Slide 6)

The “Identify Parameters and Statistics” slide was missing context that makes it clear whether each quantity describes a population or a sample. The updated examples are:

1. Western University has 30,000 students total. The registrar has complete enrollment records for every one of them. **The average height of all 30,000 Western students is 170 cm.**
→ Parameter
2. Western has 30,000 students, but you only have time to study a small group. **You measure the height of 50 randomly selected students and find their average height is 172 cm.**
→ Statistic
3. Every person on Earth is either left-handed or not. Researchers have catalogued handedness across the entire world population. **The true proportion of left-handed people worldwide is 10%.**
→ Parameter
4. There are millions of eligible voters in the country, but only a small fraction were contacted. **In a poll of 500 voters, 52% support a policy.**
→ Statistic

Notation: Parameters vs. Statistics

Some examples of common parameters and statistics of interest are:

Measure	Population (Parameter)	Sample (Statistic)
Mean	μ	\bar{x}
Standard Deviation	σ	s
Proportion	p	$\hat{p} = \frac{\sum_{i=1}^n x_i}{n}$

From chapter 2

Note: \bar{x}_n - the observed sample mean for a sample of size n
e.g. we might see $\bar{x}_3 = 172.875$
 $\bar{x}^{(k)}$ - the observed sample mean for some prespecified sample size n .

Example of a probability of interest being approximated:

- Suppose we toss a coin and let X record the result as follows:

$$X_i = \begin{cases} 1, & \text{if } i^{\text{th}} \text{ toss is Heads} \\ 0, & \text{if } i^{\text{th}} \text{ toss is Tails} \end{cases}$$

Then $p := P(X_i = 1) = P(\text{Heads})$ is a parameter (true proportion/prob^y of heads)

- Given a sequence of coin tosses X_1, \dots, X_n , where X_i is the result of the i^{th} coin toss, $\hat{p} = \frac{\sum_{i=1}^n X_i}{n}$ is our sample proportion.

Statistic

Focus: the sample mean as an estimator of the population mean

We often use the sample mean \bar{X}_n to estimate the population mean μ .

But is this a good idea?

PART 2

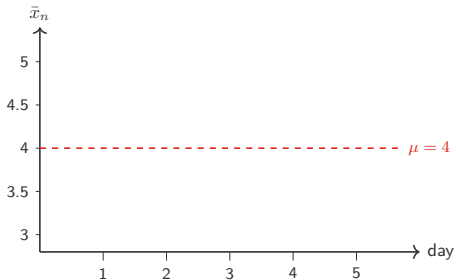
Does \bar{X}_n Aim at the Right Target?

The Law of Large Numbers

Building the Running Mean

Example 15.2: Days 1–5

Context: A café claims the average wait is $\mu = 4$ min. You time your wait each morning.



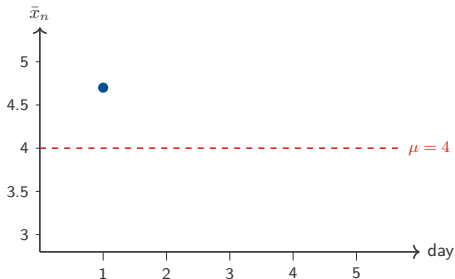
Building the Running Mean

Example 15.2: Days 1–5

Context: A café claims the average wait is $\mu = 4$ min. You time your wait each morning.

Day 1: waited 4.7 min

$$\bar{x}_1 = 4.7/1 = 4.70$$



Building the Running Mean

Example 15.2: Days 1–5

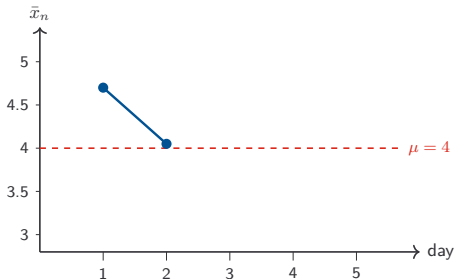
Context: A café claims the average wait is $\mu = 4$ min. You time your wait each morning.

Day 1: waited 4.7 min

$$\bar{x}_1 = 4.7/1 = 4.70$$

Day 2: waited 3.4 min

$$\bar{x}_2 = 8.1/2 = 4.05$$



Building the Running Mean

Example 15.2: Days 1–5

Context: A café claims the average wait is $\mu = 4$ min. You time your wait each morning.

Day 1: waited 4.7 min

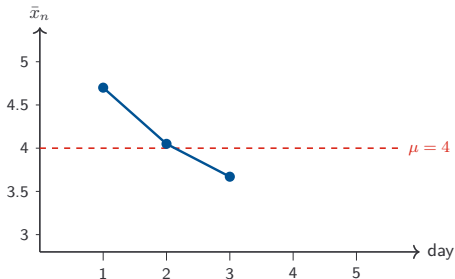
$$\bar{x}_1 = 4.7/1 = 4.70$$

Day 2: waited 3.4 min

$$\bar{x}_2 = 8.1/2 = 4.05$$

Day 3: waited 2.9 min

$$\bar{x}_3 = 11.0/3 = 3.67$$



Building the Running Mean

Example 15.2: Days 1–5

Context: A café claims the average wait is $\mu = 4$ min. You time your wait each morning.

Day 1: waited 4.7 min

$$\bar{x}_1 = 4.7/1 = 4.70$$

Day 2: waited 3.4 min

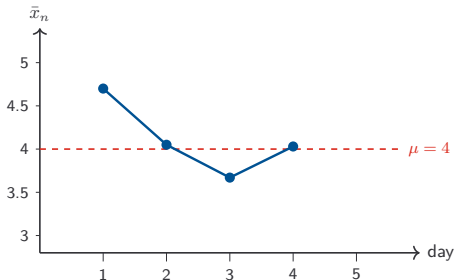
$$\bar{x}_2 = 8.1/2 = 4.05$$

Day 3: waited 2.9 min

$$\bar{x}_3 = 11.0/3 = 3.67$$

Day 4: waited 5.1 min

$$\bar{x}_4 = 16.1/4 = 4.03$$



Building the Running Mean

Example 15.2: Days 1–5

Context: A café claims the average wait is $\mu = 4$ min. You time your wait each morning.

Day 1: waited 4.7 min

$$\bar{x}_1 = 4.7/1 = 4.70$$

Day 2: waited 3.4 min

$$\bar{x}_2 = 8.1/2 = 4.05$$

Day 3: waited 2.9 min

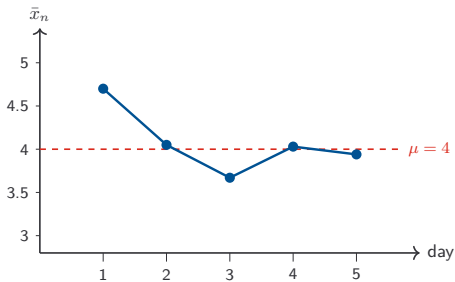
$$\bar{x}_3 = 11.0/3 = 3.67$$

Day 4: waited 5.1 min

$$\bar{x}_4 = 16.1/4 = 4.03$$

Day 5: waited 3.6 min

$$\bar{x}_5 = 19.7/5 = 3.94$$

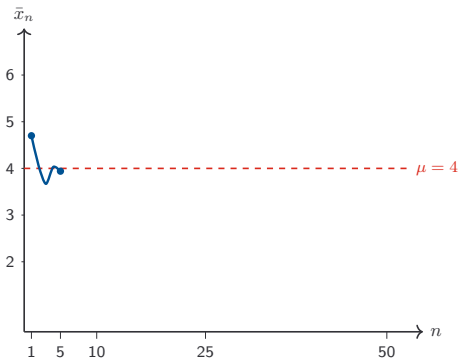


The Law of Large Numbers in Action

Example 15.2: Convergence

Context: After 5 mornings, $\bar{x}_5 = 3.94$.

Keep going...



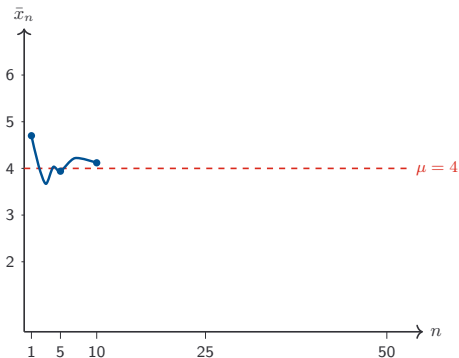
The Law of Large Numbers in Action

Example 15.2: Convergence

Context: After 5 mornings, $\bar{x}_5 = 3.94$.

Keep going...

Day 10: $\bar{x}_{10} = 4.12$ (closer)



The Law of Large Numbers in Action

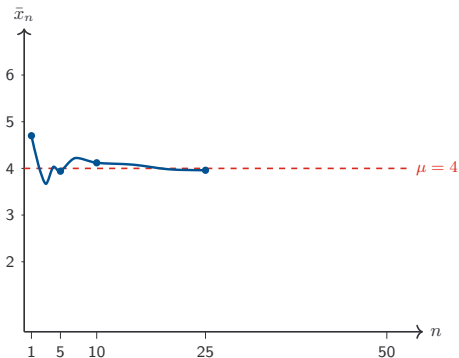
Example 15.2: Convergence

Context: After 5 mornings, $\bar{x}_5 = 3.94$.

Keep going...

Day 10: $\bar{x}_{10} = 4.12$ (closer)

Day 25: $\bar{x}_{25} = 3.96$ (very close)



The Law of Large Numbers in Action

Example 15.2: Convergence

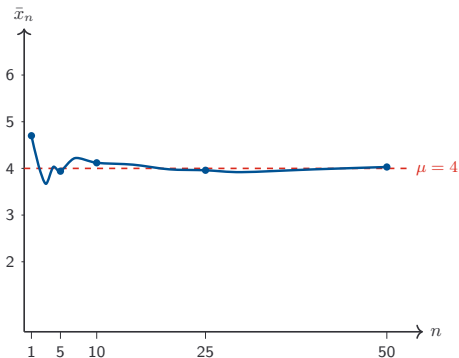
Context: After 5 mornings, $\bar{x}_5 = 3.94$.

Keep going...

Day 10: $\bar{x}_{10} = 4.12$ (closer)

Day 25: $\bar{x}_{25} = 3.96$ (very close)

Day 50: $\bar{x}_{50} = 4.03$ (essentially at μ)



The Law of Large Numbers in Action

Example 15.2: Convergence

Context: After 5 mornings, $\bar{x}_5 = 3.94$.

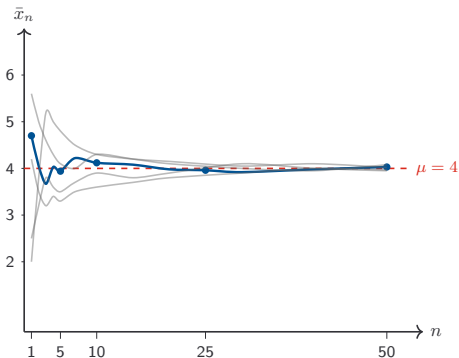
Keep going...

Day 10: $\bar{x}_{10} = 4.12$ (closer)

Day 25: $\bar{x}_{25} = 3.96$ (very close)

Day 50: $\bar{x}_{50} = 4.03$ (essentially at μ)

This happens **every time** you repeat the experiment.



The Law of Large Numbers

Law of Large Numbers (LLN)

If X_1, X_2, \dots, X_n is a random sample from a population with mean μ , then as the sample size n increases, the sample mean \bar{X}_n converges to μ with high probability.

In plain language:

- Larger samples give **better estimates**

The Law of Large Numbers

Law of Large Numbers (LLN)

If X_1, X_2, \dots, X_n is a random sample from a population with mean μ , then as the sample size n increases, the sample mean \bar{X}_n converges to μ with high probability.

In plain language:

- Larger samples give **better estimates**
- \bar{X}_n gets closer to μ as n grows


The Law of Large Numbers

Law of Large Numbers (LLN)

If X_1, X_2, \dots, X_n is a random sample from a population with mean μ , then as the sample size n increases, the sample mean \bar{X}_n converges to μ with high probability.

In plain language:

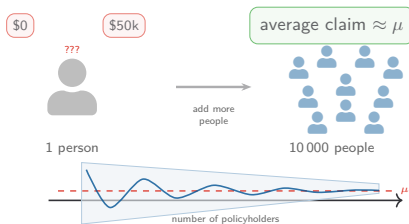
- Larger samples give **better estimates**
- \bar{X}_n gets closer to μ as n grows
- This justifies using \bar{X}_n to estimate μ

 **Intuition:** The LLN tells us *where* \bar{X}_n converges (to μ), but not *how much* it varies along the way. To understand variability, we need its **sampling distribution**.

LLN in Practice: Insurance & Risk Pooling

The core idea:

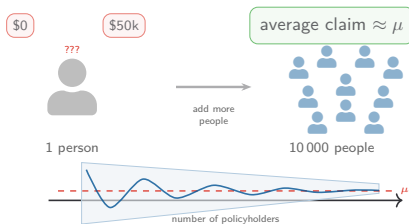
- Cannot predict whether *you* will file a claim
- With thousands of policyholders, the *average* claim cost converges reliably to μ



LLN in Practice: Insurance & Risk Pooling

The core idea:

- Cannot predict whether *you* will file a claim
- With thousands of policyholders, the *average* claim cost converges reliably to μ



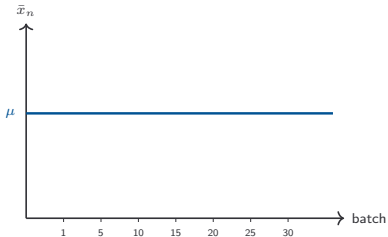
💡 Intuition: This is why insurance *works*: the company doesn't need to predict your outcome - only the average across the pool.

LLN in Practice: Quality Control

The core idea:

- A factory cannot test every unit

Control Chart (Schematic)

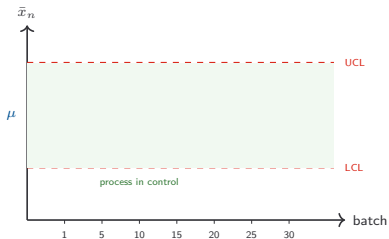


LLN in Practice: Quality Control

The core idea:

- A factory cannot test every unit
- Measures the average defect rate across each production run

Control Chart (Schematic)



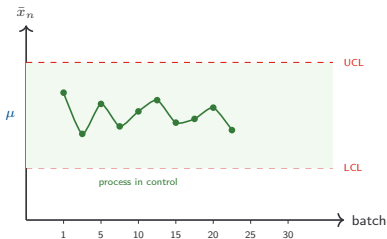
Generally, control limits are set at $\mu \pm 3 \cdot SE$.

LLN in Practice: Quality Control

The core idea:

- A factory cannot test every unit
- Measures the average defect rate across each production run
- As the run grows, the sample average tracks the true defect rate, catching process drift early

Control Chart (Schematic)



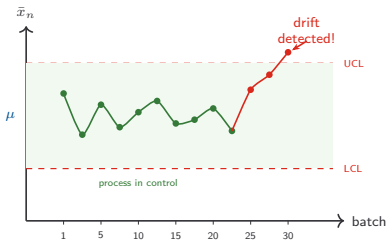
Generally, control limits are set at $\mu \pm 3 \cdot SE$.

LLN in Practice: Quality Control

The core idea:

- A factory cannot test every unit
- Measures the average defect rate across each production run
- As the run grows, the sample average tracks the true defect rate, catching process drift early

Control Chart (Schematic)



Generally, control limits are set at $\mu \pm 3 \cdot SE$.

💡 Intuition: If the running average suddenly deviates from the known μ , something has *changed* in the production process.

This is the basis of **control charts** in statistical process control.

What the LLN Does *Not* Say

Common Misconception: The Gambler's Fallacy

"I've flipped 10 tails in a row. Heads must be due next."

What the LLN actually says:

As $n \rightarrow \infty$, $\bar{X}_n \rightarrow \mu$ because early deviations get *diluted*, not corrected.

Each flip is independent: past outcomes carry no information about future ones.

Real-world stakes:

- Casinos profit because gamblers misread the LLN as a short-run correction mechanism.
- Investors who "average down" after repeated losses may be committing the same error.
- The LLN is a *long-run* guarantee, not a short-run balancing force.

PART 3

How Much Does \bar{X}_n Vary?

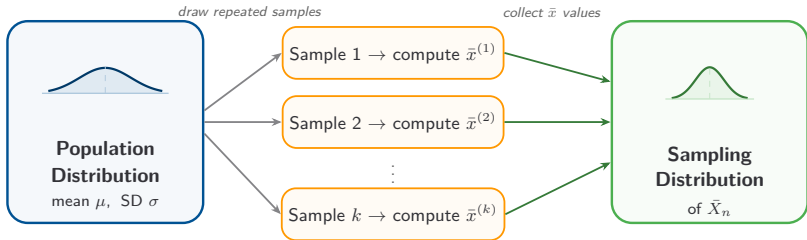
Sampling distributions and the standard error

What Is a Sampling Distribution?

Sampling Distribution

↙ sample mean

The **sampling distribution** of a statistic is the distribution of values taken by the statistic in *all possible samples* of the same size from the same population.



In theory, k should be the size of all possible samples from the popⁿ.

Sampling Distribution by Hand ($n = 2$)

Context: A fair coin is tossed. Let $X = 0$ (tails) or $X = 1$ (heads), each with probability $1/2$.

All $2^2 = 4$ possible samples of size 2:

Sample	\bar{x}
(0, 0)	0
(0, 1)	0.5
(1, 0)	0.5
(1, 1)	1

Sampling Distribution by Hand ($n = 2$)

Context: A fair coin is tossed. Let $X = 0$ (tails) or $X = 1$ (heads), each with probability $1/2$.

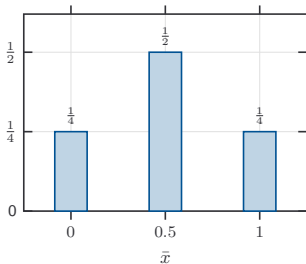
All $2^2 = 4$ possible samples of size 2:

Sample	\bar{x}	Probability
(0, 0)	0	$1/4 = 0.25$
(0, 1)	0.5	$1/4 = 0.25$
(1, 0)	0.5	$1/4 = 0.25$
(1, 1)	1	$1/4 = 0.25$

$\bar{x} = 0.5$ ←

0.5

Sampling distribution of \bar{X}_2 :



Sampling Distribution by Hand ($n = 3$)

Context: Same coin flip: $X = 0$ or 1 , each with probability $1/2$.

All $2^3 = 8$ possible samples of size 3:

Sample	\bar{x}
(0, 0, 0)	0
(0, 0, 1)	$1/3$
(0, 1, 0)	$1/3$
(1, 0, 0)	$1/3$
(0, 1, 1)	$2/3$
(1, 0, 1)	$2/3$
(1, 1, 0)	$2/3$
(1, 1, 1)	1

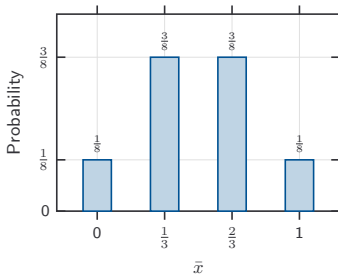
Sampling Distribution by Hand ($n = 3$)

Context: Same coin flip: $X = 0$ or 1 , each with probability $1/2$.

All $2^3 = 8$ possible samples of size 3:

Sample	\bar{x}
(0, 0, 0)	0
(0, 0, 1)	$1/3$
(0, 1, 0)	$1/3$
(1, 0, 0)	$1/3$
(0, 1, 1)	$2/3$
(1, 0, 1)	$2/3$
(1, 1, 0)	$2/3$
(1, 1, 1)	1

Sampling distribution of \bar{X}_3 :

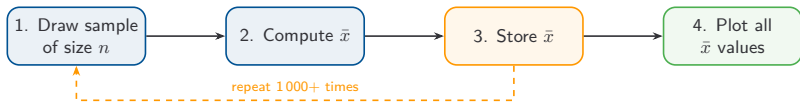


From Enumeration to Simulation

With only a few outcomes, we can list every possible sample. For larger or continuous populations, we **simulate** instead:

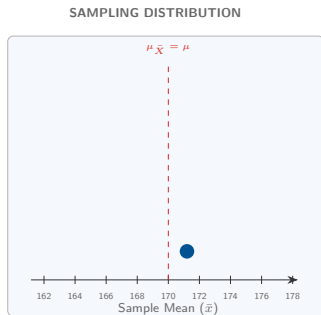
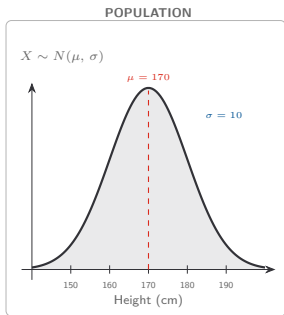
From Enumeration to Simulation

With only a few outcomes, we can list every possible sample. For larger or continuous populations, we **simulate** instead:

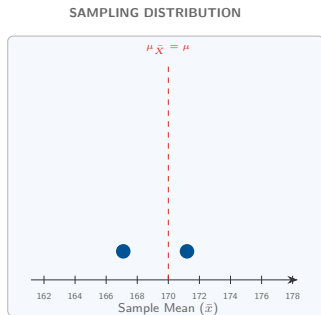
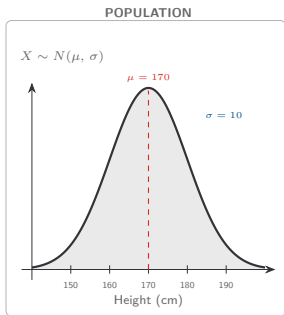


🗨 **Note:** The result is technically an **empirical** sampling distribution. With enough repetitions, it closely approximates the true (theoretical) sampling distribution.

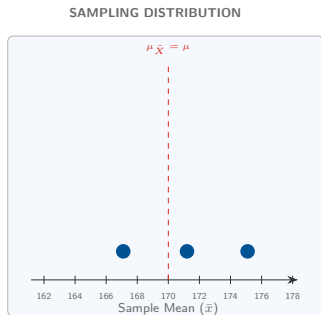
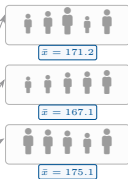
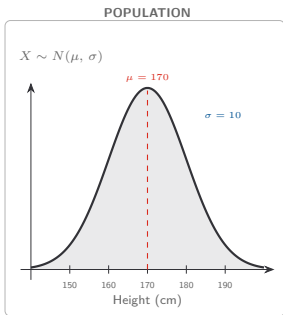
Building a Sampling Distribution



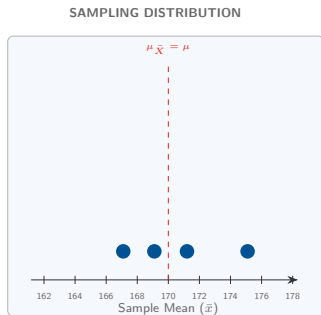
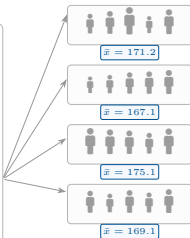
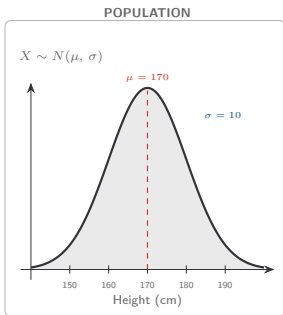
Building a Sampling Distribution



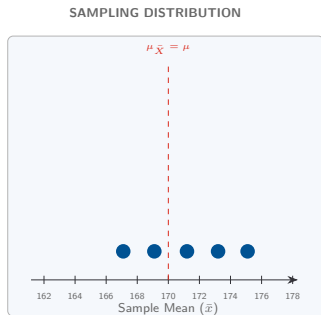
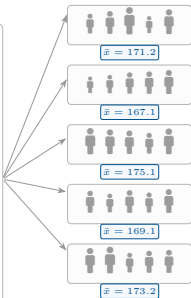
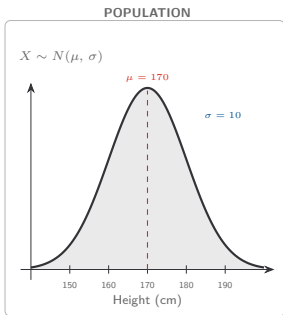
Building a Sampling Distribution



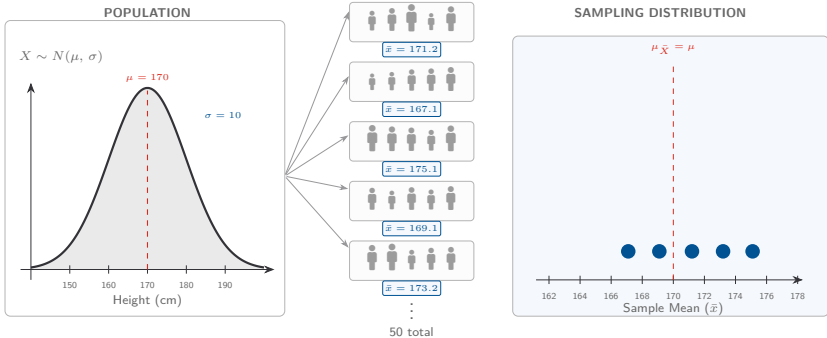
Building a Sampling Distribution



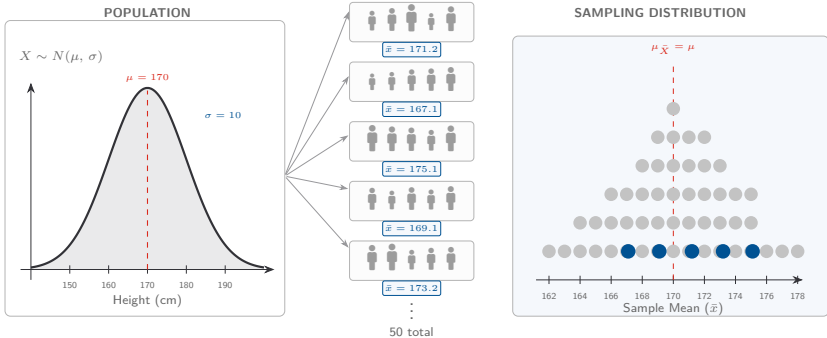
Building a Sampling Distribution



Building a Sampling Distribution



Building a Sampling Distribution



How does the sample mean vary?

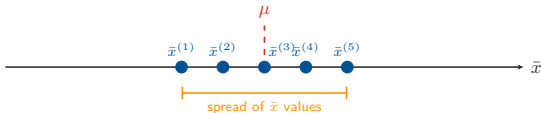
Why Sampling Distributions Matter

The LLN tells us \bar{X}_n heads toward μ as n grows, but it doesn't tell us how *close* we can expect \bar{X}_n to be for a given sample size.

Why Sampling Distributions Matter

The LLN tells us \bar{X}_n heads toward μ as n grows, but it doesn't tell us how *close* we can expect \bar{X}_n to be for a given sample size.

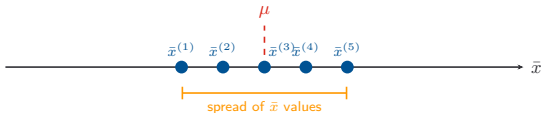
Imagine drawing five different samples of the same size from the same population. Each gives a slightly different \bar{x} :



Why Sampling Distributions Matter

The LLN tells us \bar{X}_n heads toward μ as n grows, but it doesn't tell us how *close* we can expect \bar{X}_n to be for a given sample size.

Imagine drawing five different samples of the same size from the same population. Each gives a slightly different \bar{x} :



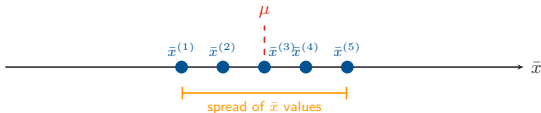
Two natural questions:

- How spread out are these sample means?

Why Sampling Distributions Matter

The LLN tells us \bar{X}_n heads toward μ as n grows, but it doesn't tell us how *close* we can expect \bar{X}_n to be for a given sample size.

Imagine drawing five different samples of the same size from the same population. Each gives a slightly different \bar{x} :

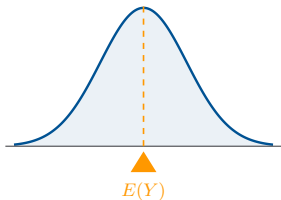


Two natural questions:

- How spread out are these sample means?
- What is the relationship between the spread and the sample size n ?

Recap: Centre and Spread

Centre: $E(Y)$



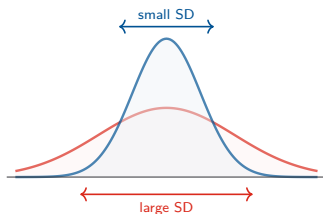
- Expected value
- Expectation
- Average of the
r.v. Y

The **expected value** $E(Y)$ is the balance point of the distribution: the long-run average if you could repeat the random process forever.

Recap: Centre and Spread

The $SD(Y)$ is the average deviation from Y to its expected value $E(Y)$.

Spread: $SD(Y)$



The **standard deviation** $SD(Y)$ measures how far values typically fall from the centre.

The **variance** is the same idea in squared units: $\text{Var}(Y) = [SD(Y)]^2$.

Standard Error vs. Standard Deviation


Standard Error (SE)

The **standard error** of a statistic is the standard deviation of its *sampling distribution*. For the sample mean:

$$SE(\bar{X}_n) = SD(\bar{X}_n) = \frac{\sigma}{\sqrt{n}}$$

Why not just say “standard deviation”?

- $\sigma = SD(X)$ describes variability of **individual observations** in the population.
- $SE = \sigma/\sqrt{n}$ describes variability of **sample means** across repeated samples.

 **Intuition:** Both are standard deviations, but of *different things*: σ measures spread of raw data; SE measures precision of a statistic.

Mean and Standard Error of \bar{X}_n

Mean and Standard Error of the Sample Mean

If X is a population with mean μ and standard deviation σ , then the sample mean \bar{X}_n has:

- Mean:

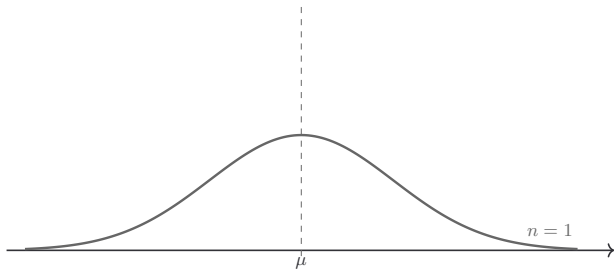
$$E(\bar{X}_n) = \mu$$

- Standard Error (SE):

$$SE(\bar{X}_n) = \sigma/\sqrt{n}$$

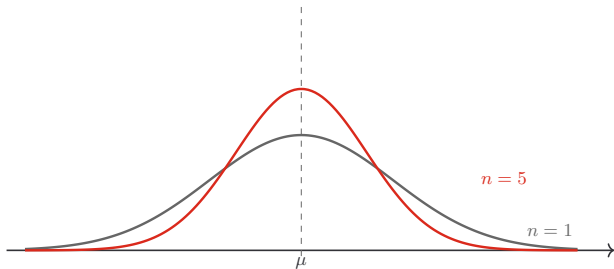
Principle: The sample mean is **unbiased** (on average equals μ), and the SE **shrinks** as n grows: larger samples give more precise estimates.

How Sample Size Affects the Sampling Distribution



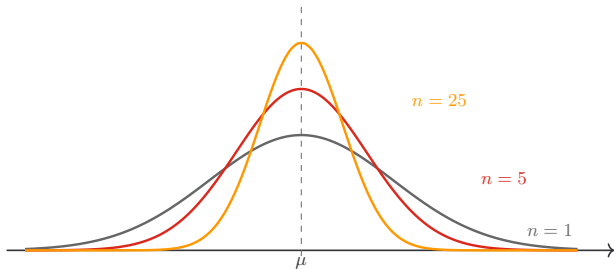
💡 Intuition: As n increases, the sampling distribution gets **narrower** (smaller SE), but the centre stays at μ .

How Sample Size Affects the Sampling Distribution



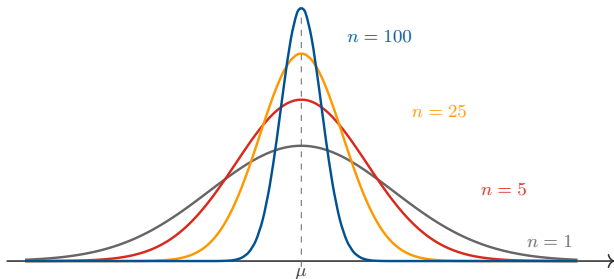
💡 Intuition: As n increases, the sampling distribution gets **narrower** (smaller SE), but the centre stays at μ .

How Sample Size Affects the Sampling Distribution



💡 Intuition: As n increases, the sampling distribution gets **narrower** (smaller SE), but the centre stays at μ .

How Sample Size Affects the Sampling Distribution



Intuition: As n increases, the sampling distribution gets **narrower** (smaller SE), but the centre stays at μ .

Finding SE's

Example 15.3

Context: Adult heights have $\mu = 170$ cm and $\sigma = 10$ cm.

Calculate: The standard error for each sample size.

Sample size n	Calculation	SE
4	$10/\sqrt{4} = 10/2$	5
9	$10/\sqrt{9} = 10/3$	3.33
25	$10/\sqrt{25} = 10/5$	2
100	$10/\sqrt{100} = 10/10$	1
400		0.5

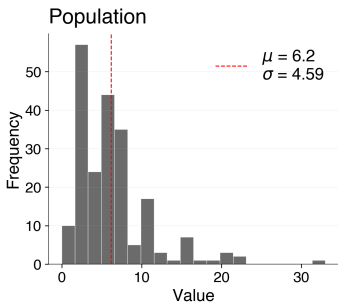
To shrink SE by a half, we need to quadruple current n

Countries Visited

Example 15.4

Context: The population distribution of countries visited by DS 1000 students ($\mu = 6.2$, $\sigma = 4.59$).

Find: The sampling distribution of \bar{X}_n for several sample sizes.



(a) What is the parameter of interest?

μ : true number of countries visited

(b) What happens to mean and SE as n increases?

$E(\bar{X}_n) = \mu = 6.2$, $SE(\bar{X}_n) = \frac{4.59}{\sqrt{n}}$
decreases as n increases.

Using the SE to Find a Minimum Sample Size

The relationship $SE = \frac{\sigma}{\sqrt{n}}$ works in both directions: given a desired standard error, we can solve for the required n .

Finding sample size achieving desired precision

Let $SE = d$

Goal: find n s.t.

$$\frac{\sigma}{\sqrt{n}} \leq d$$

$$\Rightarrow \frac{\sigma^2}{n} \leq d^2 \quad \text{by squaring}$$

$$\Rightarrow \frac{n}{\sigma^2} \leq \frac{1}{d^2} \Rightarrow \text{the smallest } n \text{ can be is } \left\lceil \frac{\sigma^2}{d^2} \right\rceil \text{ because } n \text{ is an integer.}$$

Principle: Quadrupling the sample size halves the SE: precision is gained slowly, so plan your sample size in advance.

Computing Mean and Standard Error

Example 15.5

Context: The average number of countries visited by DS 1000 students has $\mu = 6.2$ and $\sigma = 4.59$.

Find: The required sample sizes for target precision levels.

(a) What sample size is needed to halve the SE?

Let $SE_0 := \sigma / \sqrt{n_0}$ $SE_1 = \sigma / \sqrt{n_1}$,
By assumption, $SE_1 = \frac{1}{2} SE_0 \Rightarrow \frac{\sigma}{\sqrt{n_1}} = \frac{1}{2} \frac{\sigma}{\sqrt{n_0}}$
 $\Rightarrow \sqrt{n_0} \cancel{\sigma} = \sqrt{n_1} \frac{1}{2} \cancel{\sigma} \Rightarrow \sqrt{n_0} = \frac{1}{2} \sqrt{n_1}$

standard deviation of \bar{X}

(b) What sample size n is needed so that the SE is at most 0.5 countries?

$$SE \leq 0.5$$

$$\frac{\sigma}{\sqrt{n}} \leq 0.5 \Rightarrow \frac{\sigma^2}{n} \leq 0.25 \Rightarrow \frac{\sigma^2}{0.25} \leq n \quad \text{Recall that } n \text{ must integer}$$
$$\Rightarrow \left\lceil \frac{\sigma^2}{0.25} \right\rceil \leq n \Rightarrow \lceil 84.27 \rceil \leq n \Rightarrow \boxed{n=85}$$

PART 4

How is the Sample Mean Distributed?

From Normal populations to the Central Limit Theorem

Recap: What We Know So Far

So far, we answered two questions:

✓ **Does \bar{X}_n aim at the right target?**

Yes: the LLN says $\bar{X}_n \rightarrow \mu$ as n grows.

Recap: What We Know So Far

So far, we answered two questions:

✓ **Does \bar{X}_n aim at the right target?**

Yes: the LLN says $\bar{X}_n \rightarrow \mu$ as n grows.

✓ **How much does \bar{X}_n vary?**

SE = σ/\sqrt{n} . Larger $n \rightarrow$ less variability.

Recap: What We Know So Far

So far, we answered two questions:

✓ **Does \bar{X}_n aim at the right target?** Yes: the LLN says $\bar{X}_n \rightarrow \mu$ as n grows.

✓ **How much does \bar{X}_n vary?** $SE = \sigma/\sqrt{n}$. Larger $n \rightarrow$ less variability.

? **What distribution does the variability follow?**

← next

Sampling from a Normal Population

Sampling Distribution of Normal Populations

If X is Normally distributed with mean μ and standard deviation σ , then the sampling distribution of \bar{X}_n for samples of size n is:

$$\bar{X}_n \sim N(\mu, \sigma/\sqrt{n})$$

This result holds for *any* sample size $n \geq 1$.

Normal
Population

\implies
for *any* n

Normal
Sampling Distribution

Electric Vehicle Range

Example 15.6: Setup

Context: Real-world range for a popular EV model is Normally distributed with $\mu = 263$ miles and $\sigma = 25$ miles (Consumer Reports highway tests, 2023). A reviewer tests $n = 16$ vehicles.

(a) What is the sampling distribution of \bar{X}_{16} ?

$$\bar{X}_{16} \sim N(263, 25/\sqrt{16}) = N(263, 25/4) \\ = N(263, 6.25).$$

(b) What is the probability that $\bar{X}_{16} \geq 270$?

$$\begin{aligned} \text{Goal: Find } P(\bar{X}_{16} \geq 270) \\ &= P\left(\frac{\bar{X}_{16} - 263}{6.25} \geq \frac{270 - 263}{6.25}\right) \\ &= P(Z \geq 1.12) \\ &= 1 - P(Z \leq 1.12) \\ &= 1 - \Phi(1.12) \\ &= 1 - 0.8686 = 0.13 \end{aligned}$$

Electric Vehicle Range

Example 15.6: Calculation

(c) What is $P(250 < \bar{X}_{16} < 276)$?

$$\begin{aligned} \text{Goal: Find } & P(250 < \bar{X}_{16} < 276) \\ = & P\left(\frac{250-263}{6.25} < \frac{\bar{X}_{16}-263}{6.25} < \frac{276-263}{6.25}\right) \\ = & P(-2.08 < Z < 2.08) \\ = & 0.96 \end{aligned}$$

(d) The reviewer finds $\bar{x} = 280$ miles. What is the probability of observing a sample mean this high or higher under the stated specs?

$$\begin{aligned} \text{Goal: Find } & P(\bar{X}_{16} \geq 280) \\ = & P(Z \geq 2.72) \\ = & 0.0033 \end{aligned}$$

EV Range: What Sample Size? Recall: $SE = \frac{\sigma}{\sqrt{n}}$ (the standard deviation of \bar{X}_n)

Example 15.7

Context: EV range has $\mu = 263$ miles and $\sigma = 25$ miles (same as the previous example).

Find: The minimum sample size and the SE for a given n .

- (a) A reviewer wants the SE to be at most 5 miles. What is the minimum sample size?

Goal: Find n s.t. $SE \leq 5$.

Recall that $\bar{X}_n \sim N(\mu, \sigma/\sqrt{n}) = N(263, 25/\sqrt{n})$

$$SE = \frac{25}{\sqrt{n}} \leq 5 \Rightarrow 25 \leq 5\sqrt{n}$$

$$\Rightarrow 5 \leq \sqrt{n}$$

$\Rightarrow 25 \leq n$ as $\lceil 25 \rceil = 25$, 25 vehicles are sufficient.

- (b) If the reviewer tests $n = 100$ vehicles, what is the SE?

$$SE = \frac{\sigma}{\sqrt{n}} \Rightarrow SE = \frac{25}{\sqrt{100}} = \frac{25}{10} = 2.5$$

PART 5

The Central Limit Theorem

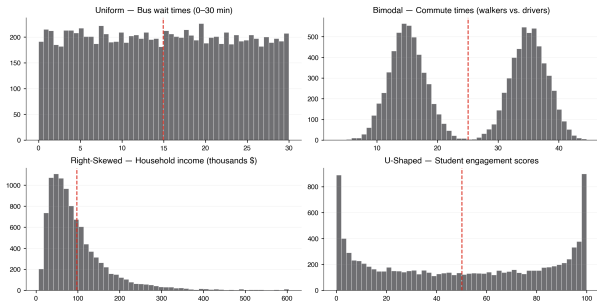
Why sample means become approximately Normal regardless of the population's distribution

Four Strange Populations

Example 15.8

None of these population distributions are normal. What will the sampling distribution of \bar{X}_n look like for large n ?

Let's investigate.



The Central Limit Theorem

Central Limit Theorem (CLT)

For a random sample of size n from *any* population with mean μ and standard deviation σ , as n increases, the sampling distribution of \bar{X}_n approaches a Normal distribution:

$$\bar{X}_n \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

where \sim means “is approximately distributed as.” The approximation improves as n increases.

Any
Population



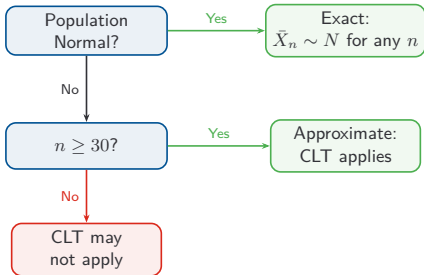
for large n

Approximately Normal
Sampling Distribution

When Does the CLT Apply?

Conditions:

1. The sample is **random**
2. The sample size is **large enough**:
 - $n \geq 30$ (rule of thumb), or
 - The population is approximately Normal



Principle: The CLT lets us use Normal probability calculations for *any* population, as long as n is large enough. If the population is already Normal, the result is exact for any n .

Potato Crate Quality Control

Example 15.9: Setup

Context: A farm ships crates of potatoes. Potato weight has $\mu = 300$ g and $\sigma = 50$ g. Each crate contains $n = 100$ potatoes. An inspector samples crates and calculates the average weight per potato \bar{X}_{100} .

(a) What are the mean and standard ~~error~~^{error} deviation of \bar{X}_{100} ?

$$E(\bar{X}_{100}) = 300, \quad SE(\bar{X}_{100}) = 50/\sqrt{100} = 50/10 = 5$$

(b) Approximate the probability that the average weight per potato in a crate is between 295 g and 305 g.

By CLT, $\bar{X}_{100} \sim N(300, 5)$
Goal: Find $P(295 \leq \bar{X}_{100} \leq 305)$
 \vdots
 $\doteq 0.67$

Potato Crate Quality Control

Example 15.9: Continued

Context: Potato weight: $\mu = 300$ g, $\sigma = 50$ g, $n = 100$; $\bar{X}_{100} \sim N(300, 5)$.

(c) Find the probability that the *total* weight of potatoes in a crate exceeds 31 kg.

• Let $T =$ total weight of potatoes in a crate

• Goal: Find $P(T \geq 31,000)$

$$T = \sum_{i=1}^{100} X_i = \bar{X}_{100} \cdot 100$$

$$\Rightarrow P(T \geq 31,000) = P(\bar{X}_{100} \cdot 100 \geq 31,000)$$

$$= P(\bar{X}_{100} \geq 310)$$

\vdots

$$\hat{=} 0.0228$$

CLT and Discrete distributions

The CLT applies to *any* distribution

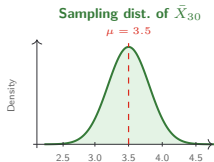
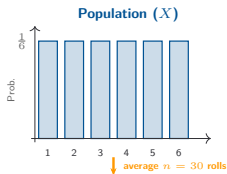
The CLT doesn't require a continuous population. Even if individual observations can only take on a **few discrete values**, the *average* of many such observations is approximately Normal.

Examples:

- Die rolls ($X \in \{1, \dots, 6\}$)
- Coin flips ($X \in \{0, 1\}$)

As long as n is large enough, \bar{X}_n is approximately Normal regardless of whether the population is discrete or continuous.

Die rolls: $\mu = 3.5$, $\sigma \approx 1.71$



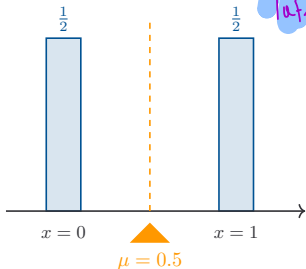
Digression: Mean of a Coin Toss

Context: Let $X = 1$ (heads) or $X = 0$ (tails), each with probability $\frac{1}{2}$.

Outcome	x	Prob.
Tails	0	$\frac{1}{2}$
Heads	1	$\frac{1}{2}$

Half the time you get 0, half the time you get 1.

The long-run average is the balance point:



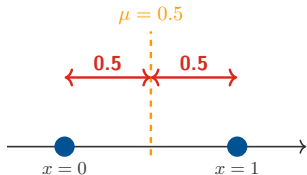
Standard Deviation of a Coin Toss

Recall: $\mu = 0.5$. How far does each flip land from the mean?

Every flip lands exactly 0.5 away from

μ :

- Tails ($x = 0$):
distance = $|0 - 0.5| = 0.5$
- Heads ($x = 1$):
distance = $|1 - 0.5| = 0.5$



The typical distance from the mean:

Principle: In general, for a binary random variable having outcomes $\{0, 1\}$ where $P(X = 1) = p$, the mean is $E(X) = \mu = p$ and the SD is $\sigma = \sqrt{p(1 - p)}$.

Coin Tosses and the CLT

Example 15.10: Setup

Context: A fair coin has $\mu = p = 0.5$ and $\sigma = 0.5$. We toss it $n = 100$ times and observe 60 heads ($\hat{p} = \bar{X}_{100} = 0.60$).

(a) What are the mean and SE of the sampling distribution of \bar{X}_{100} ?

(b) What is the sampling distribution of \bar{X}_{100} ?

Coin Tosses and the CLT

Example 15.10: Calculation

Context: Fair coin, $n = 100$; $\bar{X}_{100} \sim N(0.5, 0.05)$. Observed $\hat{p} = 0.60$.

(c) What is $P(\bar{X}_{100} \geq 0.60)$?



Putting It All Together

Example 15.11

Task: For each population and sample size, write the sampling distribution of \bar{X}_n and compute the SE.

Population	μ	σ	n	SE	Sampling dist.
Heights (Normal)	170	10	64	<input type="text"/>	<input type="text"/>
EV range (Normal)	263	25	49	<input type="text"/>	<input type="text"/>
Countries (skewed)	6.2	4.59	100	<input type="text"/>	<input type="text"/>
Income (skewed)	95.1	74.73	50	<input type="text"/>	<input type="text"/>

Principle: Normal population \Rightarrow exact (\sim).

Non-Normal population with $n \geq 30 \Rightarrow$ approximately normal.

CHAPTER 15

Summary

Key ideas and formulas to carry forward

Key Takeaways

- **Parameters vs. Statistics:** A parameter (μ, σ, p) describes a population and is fixed but unknown. A statistic (\bar{x}, s, \hat{p}) is computed from a sample and varies from sample to sample.

Key Takeaways

- **Parameters vs. Statistics:** A parameter (μ, σ, p) describes a population and is fixed but unknown. A statistic (\bar{x}, s, \hat{p}) is computed from a sample and varies from sample to sample.
- **Law of Large Numbers:** As the sample size grows, \bar{X}_n converges to μ . Larger samples produce more reliable estimates.

Key Takeaways

- **Parameters vs. Statistics:** A parameter (μ, σ, p) describes a population and is fixed but unknown. A statistic (\bar{x}, s, \hat{p}) is computed from a sample and varies from sample to sample.
- **Law of Large Numbers:** As the sample size grows, \bar{X}_n converges to μ . Larger samples produce more reliable estimates.
- **Standard Error:** The $SE = \sigma/\sqrt{n}$ measures how much \bar{X}_n varies across samples. Quadrupling n halves the SE.

Key Takeaways

- **Parameters vs. Statistics:** A parameter (μ, σ, p) describes a population and is fixed but unknown. A statistic (\bar{x}, s, \hat{p}) is computed from a sample and varies from sample to sample.
- **Law of Large Numbers:** As the sample size grows, \bar{X}_n converges to μ . Larger samples produce more reliable estimates.
- **Standard Error:** The $SE = \sigma/\sqrt{n}$ measures how much \bar{X}_n varies across samples. Quadrupling n halves the SE.
- **Sampling distribution of \bar{X}_n :**
 - If the population is Normal, then $\bar{X}_n \sim N(\mu, \sigma/\sqrt{n})$ exactly for any n .
 - If the population is not Normal, the CLT gives $\bar{X}_n \overset{\sim}{\sim} N(\mu, \sigma/\sqrt{n})$ for $n \geq 30$.

Chapter 15: Formula Card

Concept	Formula
Mean of \bar{X}_n	$E(\bar{X}_n) = \mu$
Standard Error	$SE = \frac{\sigma}{\sqrt{n}}$
Normal population	$\bar{X}_n \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$ for any n
CLT (any population)	$\bar{X}_n \dot{\sim} N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$ for $n \geq 30$
Minimum sample size	$n \geq \left(\frac{\sigma}{d}\right)^2$ where $d = \text{target SE}$

PRACTICE

Practice Questions

Ride-Share Tip Threshold

Example 15.15: Setup

Context: Tip amounts on a ride-share platform are right-skewed: $\mu = \$2.50$ and $\sigma = \$1.50$ per trip. A driver completes $n = 100$ trips per shift. The platform awards a *Top Earner* badge to drivers whose total shift tips rank in the top 2.3% under normal conditions.

(a) What are the mean and SE of \bar{X}_{100} ?

(b) Let $S = \sum_{i=1}^{100} X_i$ be the total tips for a shift. What are $E(S)$ and $SD(S)$?

Ride-Share Tip Threshold

Example 15.15: Calculation

Context: Tip amounts: $\mu = \$2.50$, $\sigma = \$1.50$, $n = 100$; $S \sim N(250, 15)$.

(c) Find the total tip threshold t^* such that $P(S > t^*) \approx 0.023$. Drivers above this earn the badge.

