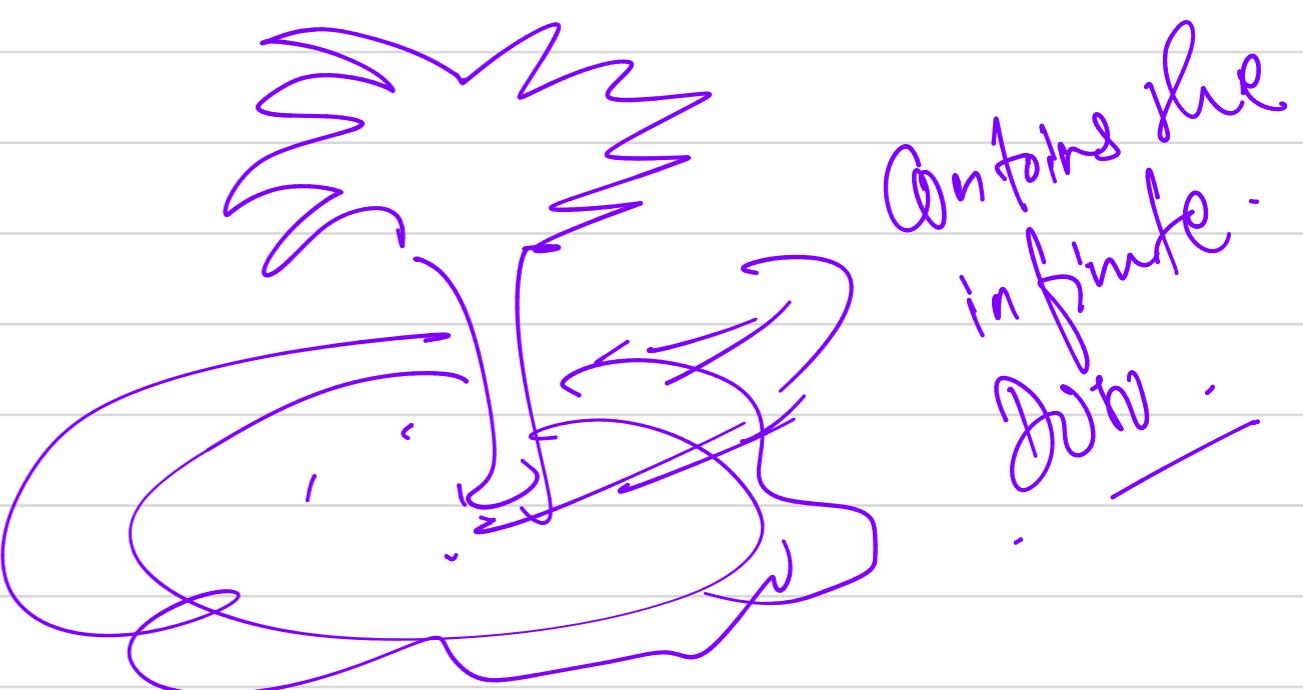


Day 79, Feb 16, 2023 (Falgun 4, 2081)

↳ Asymptotics

Asymptotics is the term for the behaviour of statistics as the sample size (or some other relevant quantity) limits to infinity (or some other relevant number).

↳ Asymptotics are incredibly useful for simple statistical inference and approximations.



↳ Asymptotics form the basis for frequency interpretation of Probabilities (the long run proportion of times an event occurs).

Limits of Random Variables

These results allow us to talk about the larger sample distribution of sample means of collection of iid observations.

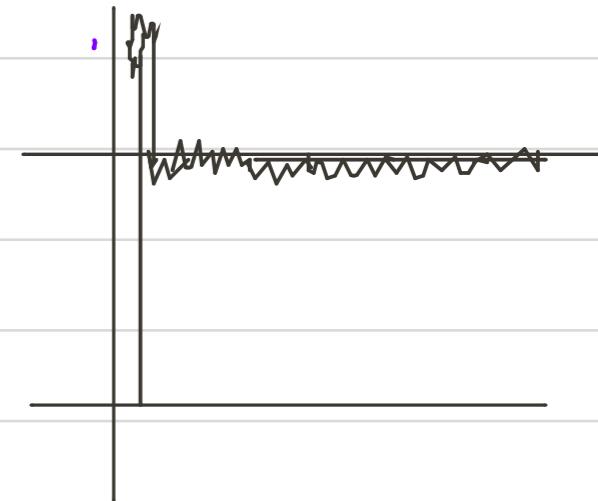
The first of these results, the Law of Large Numbers, we intuitively know

→ If soye but the average limits to what its establiy the population mean.

→ Example \bar{X}_n could be the average of the result of n coin flips (i.e. the sample proportion of heads).

→ As we flip a fair coin over and over, it eventually converges to the true probability of a head.

Law of Large Numbers in Action:



$$n \leftarrow 100$$

Mean \rightarrow means \leftarrow cumsum(rnorm(1)) / (1:n)

Ammulative sum \rightarrow generate random.

n of mean or (\bar{x}) .

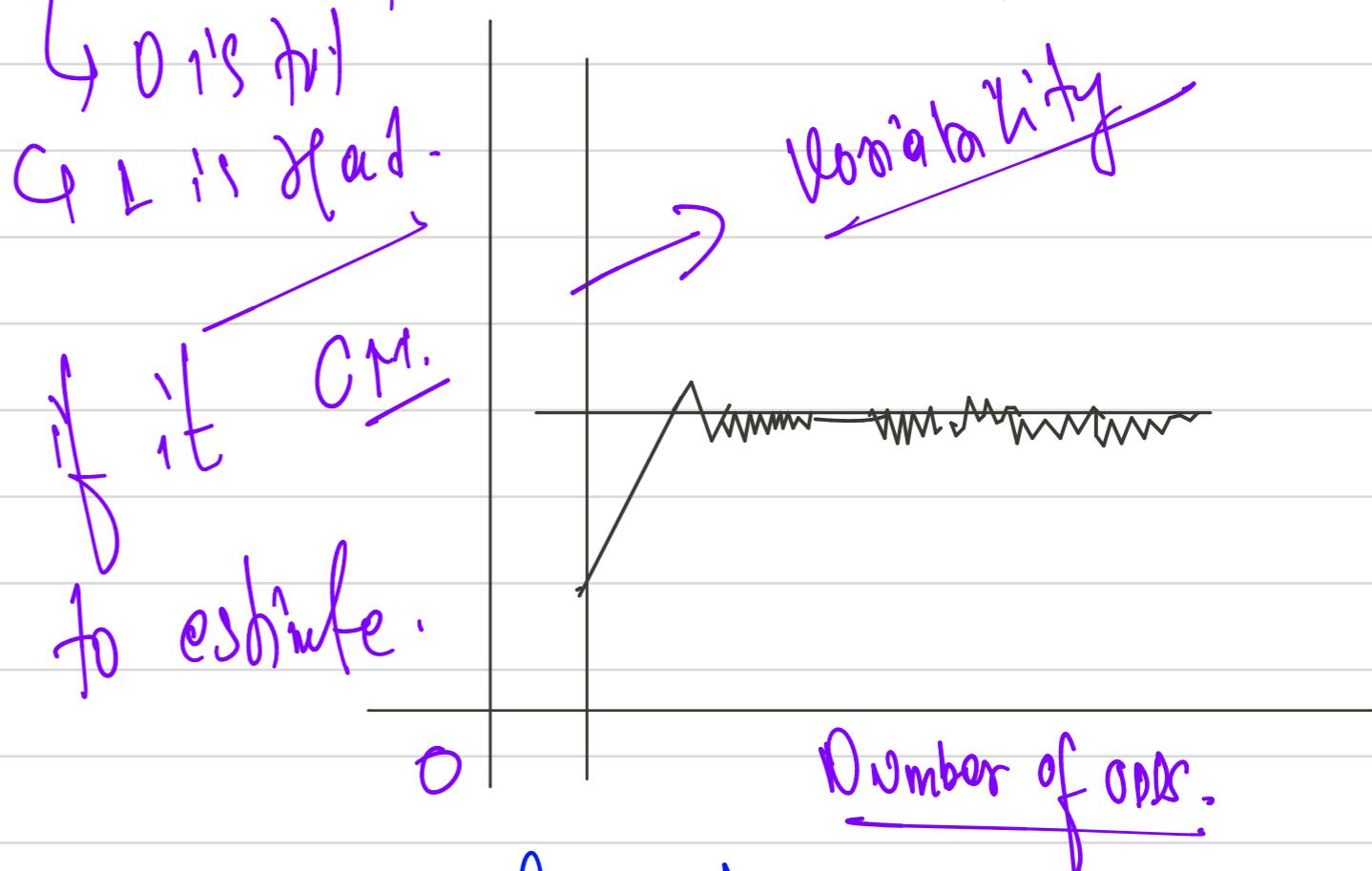
Law of Large Numbers in action, bin flip

means <- cumsum(sample(0:1, n, replace = TRUE)) / (1:n)

Discussion.

→ An estimator is consistent if it converges to what we want to estimate.

→ The LLN says that the sample mean of iid samples is consistent for the population μ .



↳ Typically, good estimators are consistent, it's not too much to ask that if we go to the trouble of collecting an infinite amount of data that we get the right answer.

⇒ The Sample Variance and the Sample Standard Deviation of iid random variables are consistent as well.

Central Limit Theorem

for our purposes, the CLT states that the distribution of averages of iid variables (properly normalized) becomes that of a standard normal as the sample size increases.

$$\frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$$\Rightarrow \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma}$$

$$\Rightarrow \frac{\text{Estimate} - \text{Mean of estimate}}{\text{Std Error of estimate}}$$

has a distribution like that of a standard normal for large n .

The useful way to think about the CLT is that \bar{X}_n is approximately $N(\mu, \sigma^2/n)$ -

Example: Simulate a Standard Normal Random Variable by rolling n (Six sided).

Let X_i be the outcome for die i

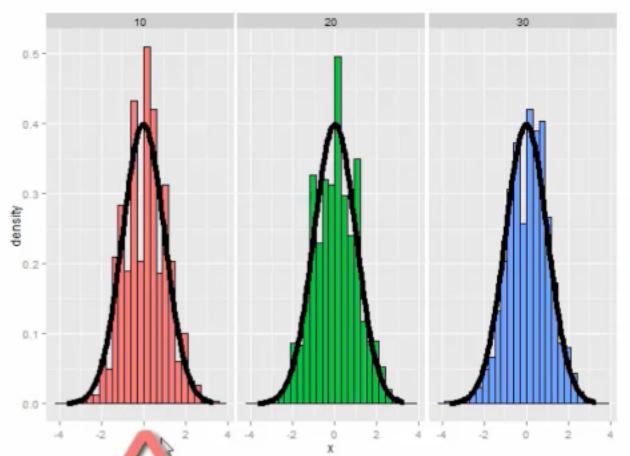
then note that $\mu = E[X_i] = 3.5$

$$\text{Var}(X_i) = 2.92$$

$$SE = \sqrt{\frac{2.92}{n}}$$

$$\Rightarrow 1.07 \pm \sqrt{\frac{2.92}{n}}$$

Result of our die rolling experiment



centered around zero because we've subtracted off the mean 3.5.

lets roll n dice, take their mean, subtract off 3.5 and divide by $\frac{1.71}{\sqrt{n}}$.

Coin CLT:

let y_i be the 0 or 1 result of the i th flip of a possibly unfair coin

- The sample proportion, say \hat{p} , is the average of the n flips

$$\rightarrow E[X_i] = p \text{ and } \text{Var}(X_i) = \frac{p(1-p)}{n}$$

$$\rightarrow \text{Standard error of the mean is } \sqrt{p(1-p)/n}$$

$$\frac{p - \bar{p}}{\sqrt{p(1-p)/n}}$$

$$\rightarrow p = \frac{1}{2}$$

$$\rightarrow p * (1-p) = \frac{1}{4}$$

$$\rightarrow \text{sqrt}(p * (1-p)) = \frac{1}{2}$$

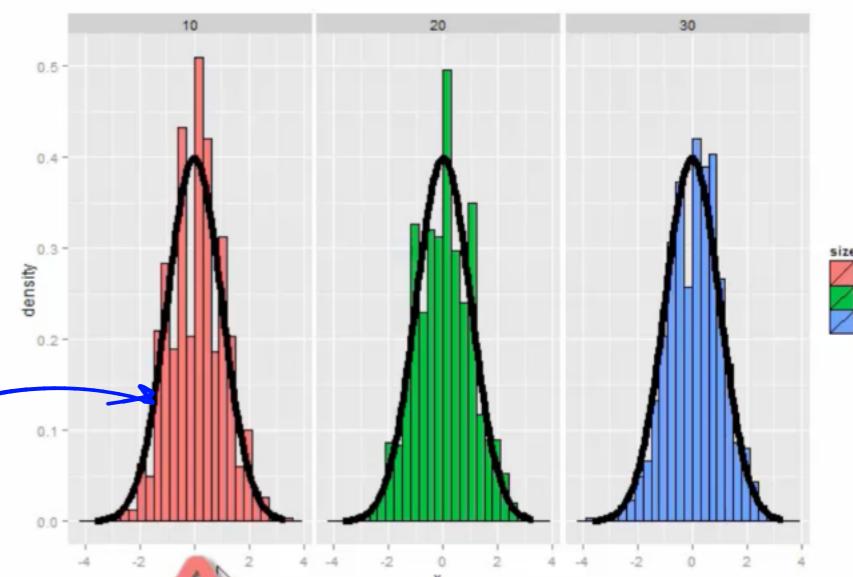
$$[p=0.9]$$

\rightarrow Simulation Results

View the

Galton's Quincunx.

Result of our die rolling experiment

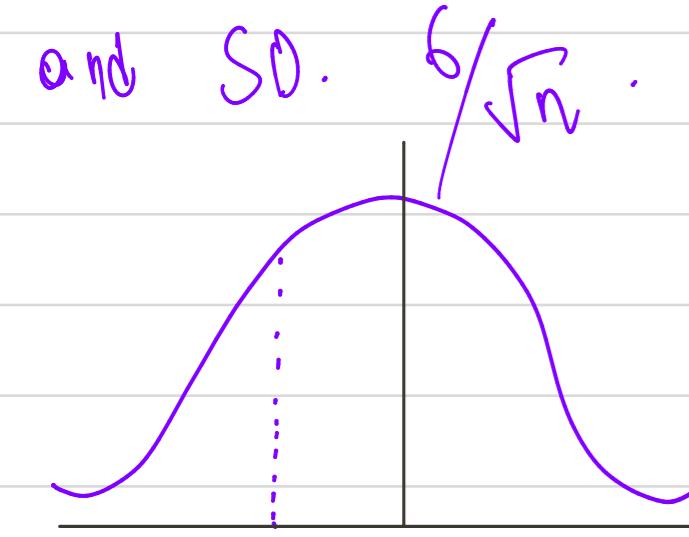


Density

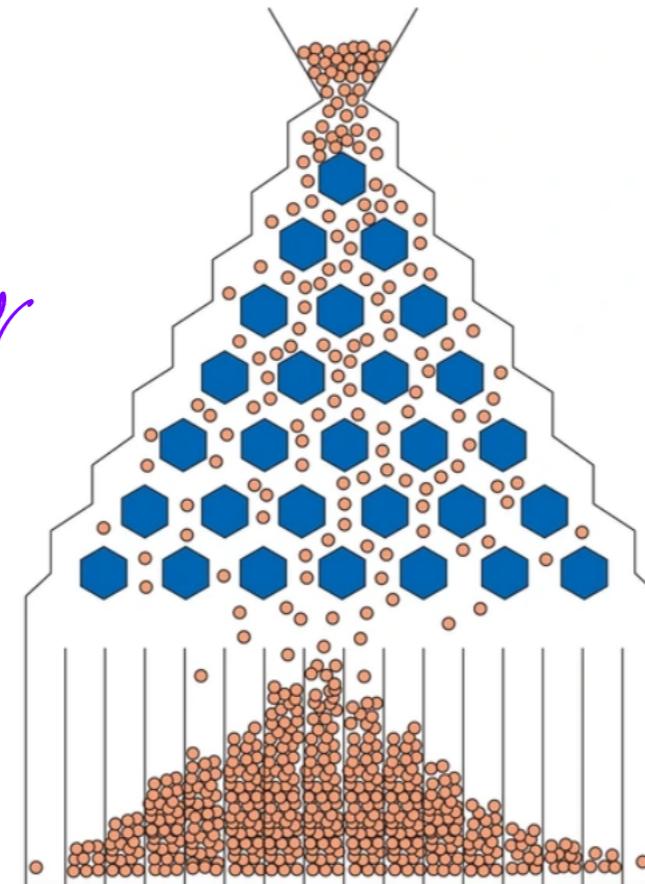
centered around zero because we've
subtracted off the mean 3.5.

Confidence Intervals

\bar{X}_n is approximately normal with μ and SD. $6/\sqrt{n}$.

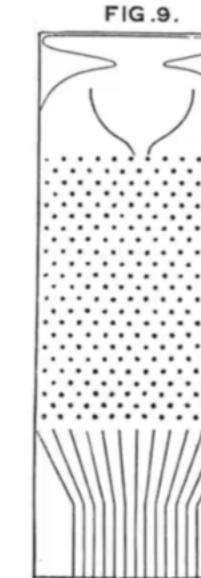
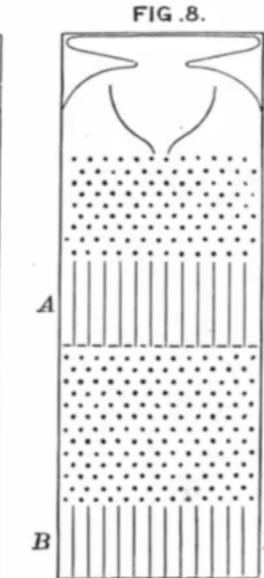


$$u - 2s / \sqrt{n}$$



→ Galton's quincunx

[http://en.wikipedia.org/wiki/Bean_machine#mediaviewer/File:Quincunx_\(Galton_Box\)_Galton_1889_diagram.png](http://en.wikipedia.org/wiki/Bean_machine#mediaviewer/File:Quincunx_(Galton_Box)_Galton_1889_diagram.png)



13/31

So, the probability \bar{X} is bigger than $\mu + 2s / \sqrt{n}$

or smaller than $\mu - 2s / \sqrt{n}$ is 5%.

$\bar{X} \pm 2s / \sqrt{n}$ is called a 95% interval for μ .

A Sample Proportions

- In the event that each X_i is 0 or 1 with common success probability p then $\sigma^2 = p(1-p)$.

- The interval takes the form

$$\hat{p} \pm z_{1-\alpha/2} \sqrt{\frac{p(1-p)}{n}}$$

- Replacing p by \hat{p} in the standard error results in what is called a Wald Confidence interval for p .

- for 95% intervals

$$\hat{p} \pm \frac{1}{\sqrt{n}}$$

is a quick CI estimate for p .

$$\sqrt{\frac{p(1-p)}{n}} \leq \sqrt{\frac{\frac{1}{2}(1-\gamma_2)}{n}} = \frac{1}{2\sqrt{n}} p(1-p)$$

$$\text{Or, } 2\sqrt{\frac{p(1-p)}{n}} \leq 2 \times \frac{1}{2\sqrt{n}} = \frac{1}{\sqrt{n}}$$

Or, for 95%.

$$p \pm \frac{1}{\sqrt{n}}$$

Simple Explanation on Confidence Intervals:

• population mean of penguin weight: 31 lbs

• population proportion of voters: 55%.

Point estimate:

↳ Uses a single value to estimate a population parameter.

Interval estimate:

↳ uses a range of values to estimate a population parameter.
↳ CI is also type of Interval Estimate.

Penguin data

• Population mean = 31 lbs .

• Sample mean = 28 lbs .

• Sample mean = 32 lbs .

Confidence Interval includes

- Sample Statistic
- Margin of Error
- Confidence Level

Interval

Sample Statistic \pm margin of error. where

Margin of error

the maximum expected difference between a population parameter and

@ Sample estimate

$$\text{Interval} : [28, 32]$$

$$30 + 2 = 32 \quad , \quad 30 - 2 = 28$$

$$\boxed{\text{Margin of Error} = Z\text{-Score} \times SE}$$

Confidence level

Describe the likelihood that a particular Sampling method will produce a Confidence Interval that includes the population parameter.

Example: 95%. CI - (95 of 100 → Contains actual population mean).

• 100 random samples

Common Confidence levels

- 90%.
- 95% → Mostly used choice -
- 99%.

Sales Revenue:

- "I think we'll do \$1,000,000 in sales."
- "Based on a 95% Confidence Interval, I estimate that our sales revenue will be between \$950,000 and \$1,050,000."

Components of a Confidence Interval

- Margin of Error:

The main components of a CI are a sample statistic, margin of error, and confidence level. Confidence Interval helps express the uncertainty of an estimate based on sample data.

- Sample Statistic and Confidence Level are other Components of CI.

Interpret Confidence Intervals

Example: Estimating the mean height of all the red maple trees in the city -

- population: 10,000 trees
- Sample: 50 trees
- Sample mean = 50 ft
- Sample Standard deviation = 1.5 ft
- Confidence Interval: 95% of CI [48, 52]

Confidence Level:

- ↳ Expresses the uncertainty of the estimation process
- ↳ 95% of intervals capture the population mean
- ↳ 5% of " do not capture the population mean.

Repeated Random Sampling is often difficult, expensive and time consuming.

↳ Confidence Intervals give data professionals a way to quantify the uncertainty due to random Sampling.

Red Maple Tree Data:

- Confidence Interval: 95% of CI [48, 52]

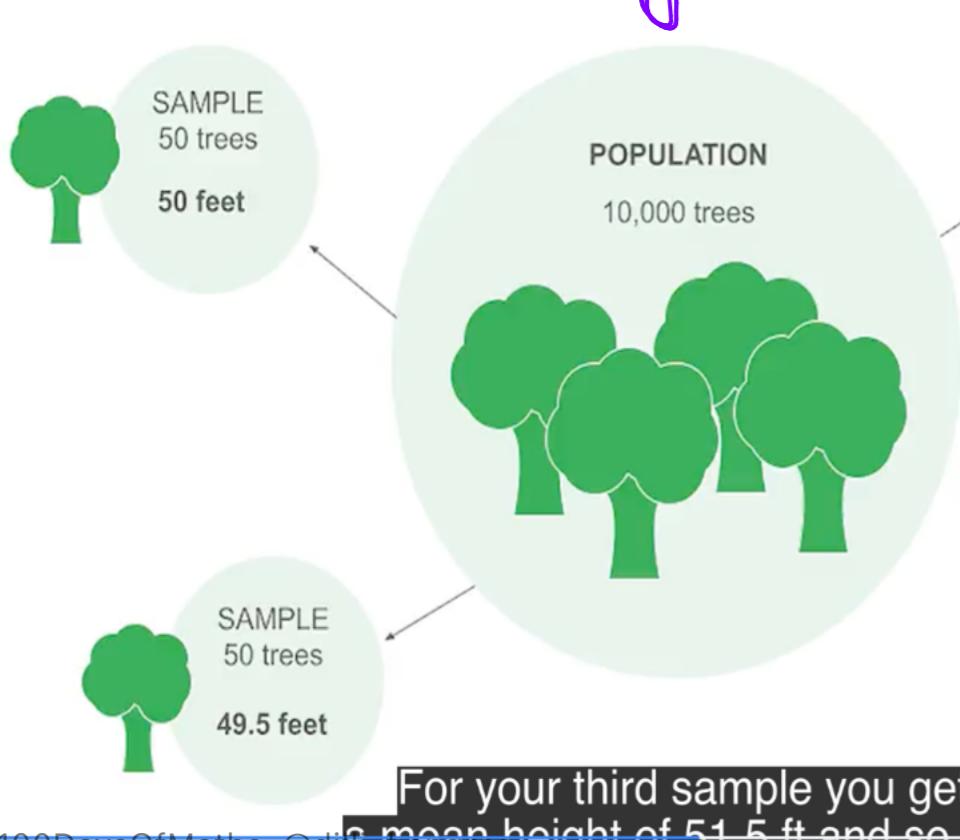
* Population mean: 52 ft.

↳ Out of 50000 trees, randomly

take 100 samples each

Sample μ values we called

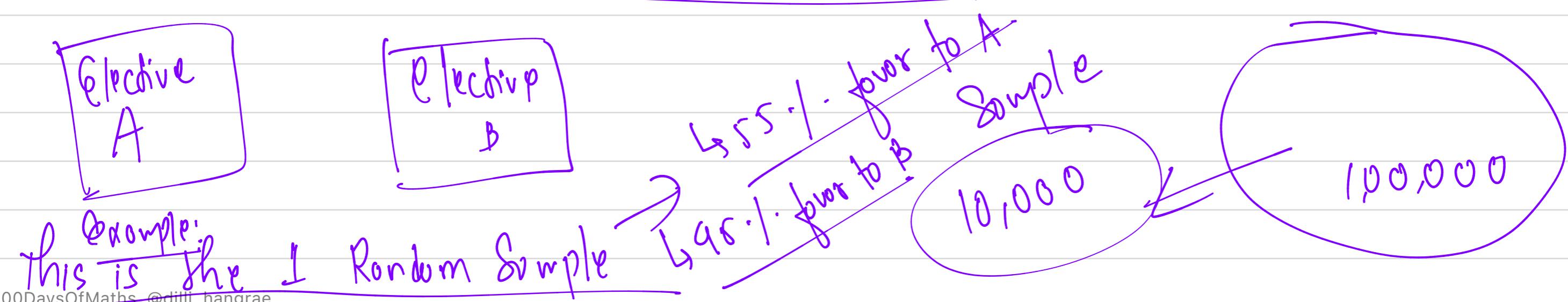
if variability:



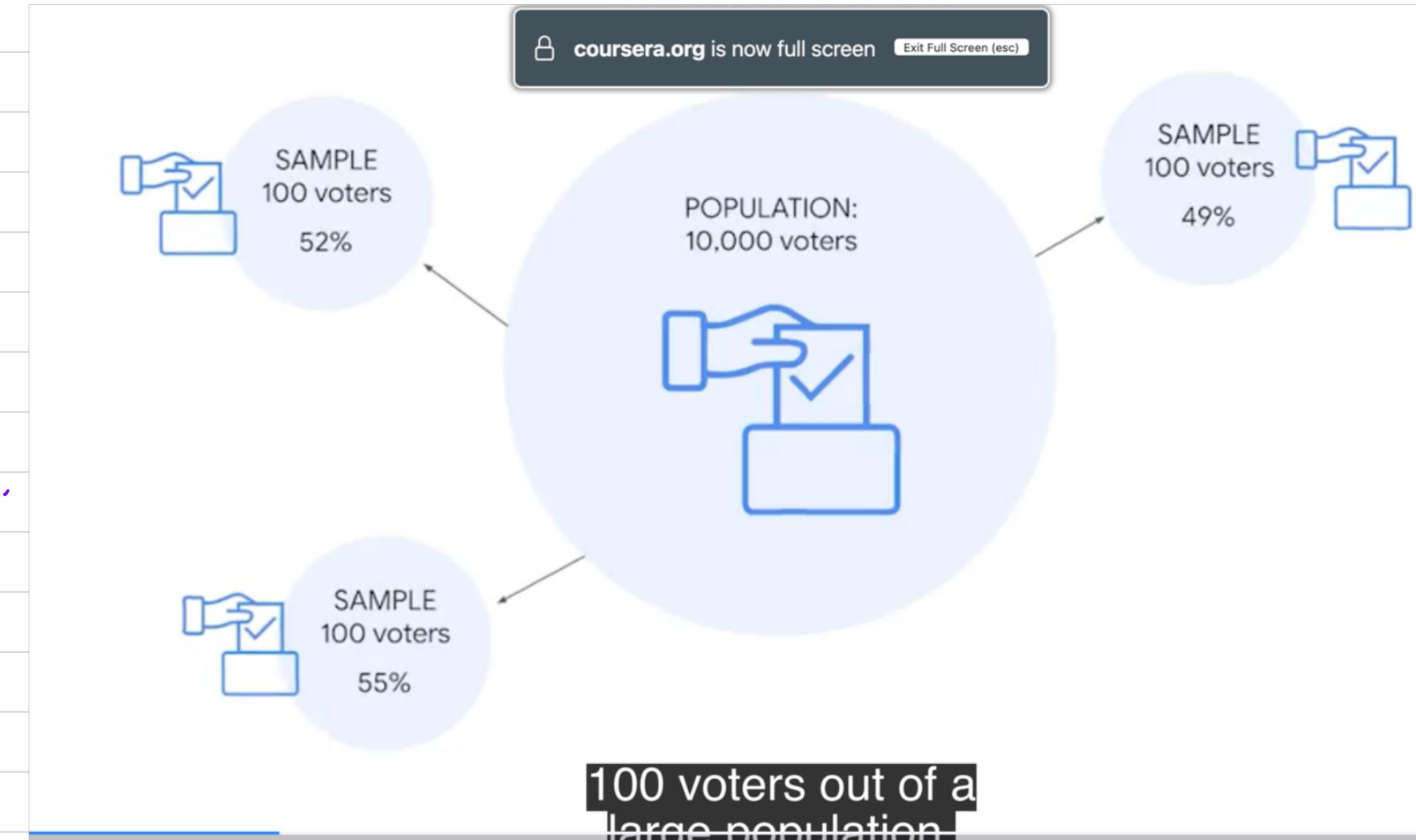
↳ Because of Sample variability any given sample will not necessarily be equal to the actual population mean.

- ① A 95% CI means that you can expect that 19 out of 20 intervals or 95% of the total will capture the population μ .
- ② 95% refer to the success rate of the process or we can expect 95% of the random intervals generate to capture the population parameter.

Construct a Confidence Interval for a Proportion



1st 100 RS
2nd 100 RS
3rd 100 RS
↓
gives different -
of favor SD
Instead of Point

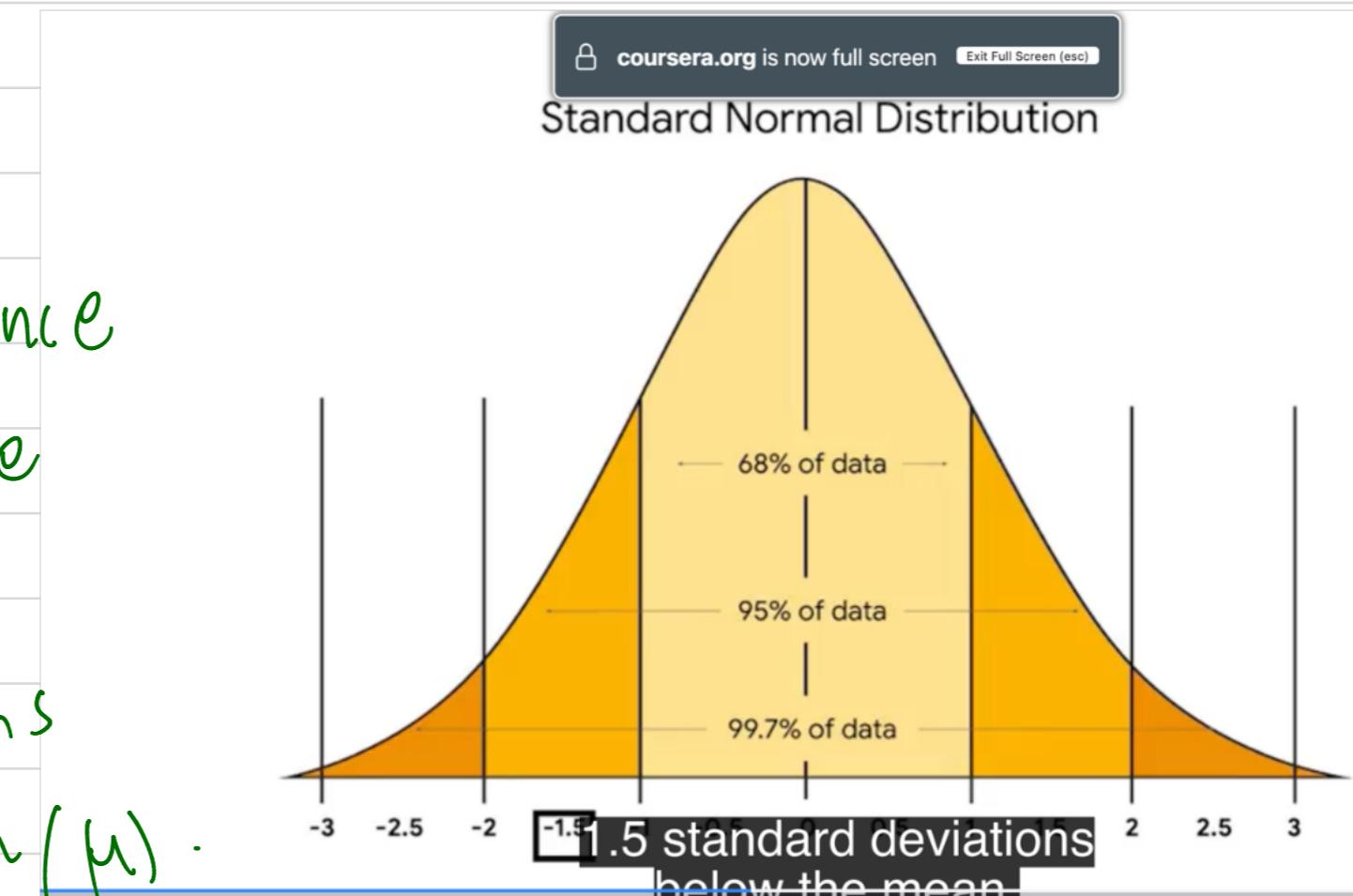


Steps for Constructing a Confidence Interval

- ① Identify a Sample Statistic : proportion
- ② Choose a Confidence Interval : 95%.
- ③ find the Margin of Error: Range of values above & below the sample statistic
- ④ Calculate the Interval:

* Margin of Error = Z-Score * SE

Z-Score measures the distance between the data point from the population mean and a standard normal distribution. Z-Score + means 1 std deviations above the mean (μ).



↳ SE measures the sample variability.

SE of the proportion

$$\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

Confidence level	Z-score
90%	1.645
95%	1.96
99%	2.58

1.645 for 90 percent

$$n = 100$$

$$\text{So, } Z\text{-Score} * \text{SE} \Rightarrow 1.96 * 0.05 \\ \Rightarrow 0.098$$

Calculate the Interval: [45.2%, 64.8%]

$$\text{Upper limit} = \text{Sample proportion} + \text{Margin of Error} \\ 0.55 + 0.098 = 0.648 \Rightarrow 64.8\%$$

Lower limit = Sample proportion - MI of Err.

$$0.55 - 0.098 \Rightarrow 0.452$$

45.2%.

$$[45.2\% < 50\%]$$

In previous elections here, the new poll reports that 54% of voters preferred candidate Davis. Using point estimation it is 50.9%. But using CI of 95%, integral will now above 50% and is 57.1%.. This gives more confidence and 5.1% of margin error need to be considered too.

$$\text{Sample Size} = 100 \quad [45.2, 64.8] \Rightarrow 19.6$$

$$\text{Sample Size} = 1000 \quad [50.9, 57.1] \Rightarrow 6.2$$

Sample Size
increases ↑

↓
st decreases

So, the first step of constructing a CI is identifying a sample
Statistic → Next, Confidence level is chosen. Then the margin of error
is found. Finally, the interval is calculated.

Construct a Confidence Interval for a Mean (μ):

↳ claim of 20 hrs of battery life is accurate for battery life of the new phone?

• Sample mean = 20.5 hrs

• Sample std. deviation = 1.7 hrs.

• Population std. deviation = 1.5 hrs.

Using Random 100 Samples

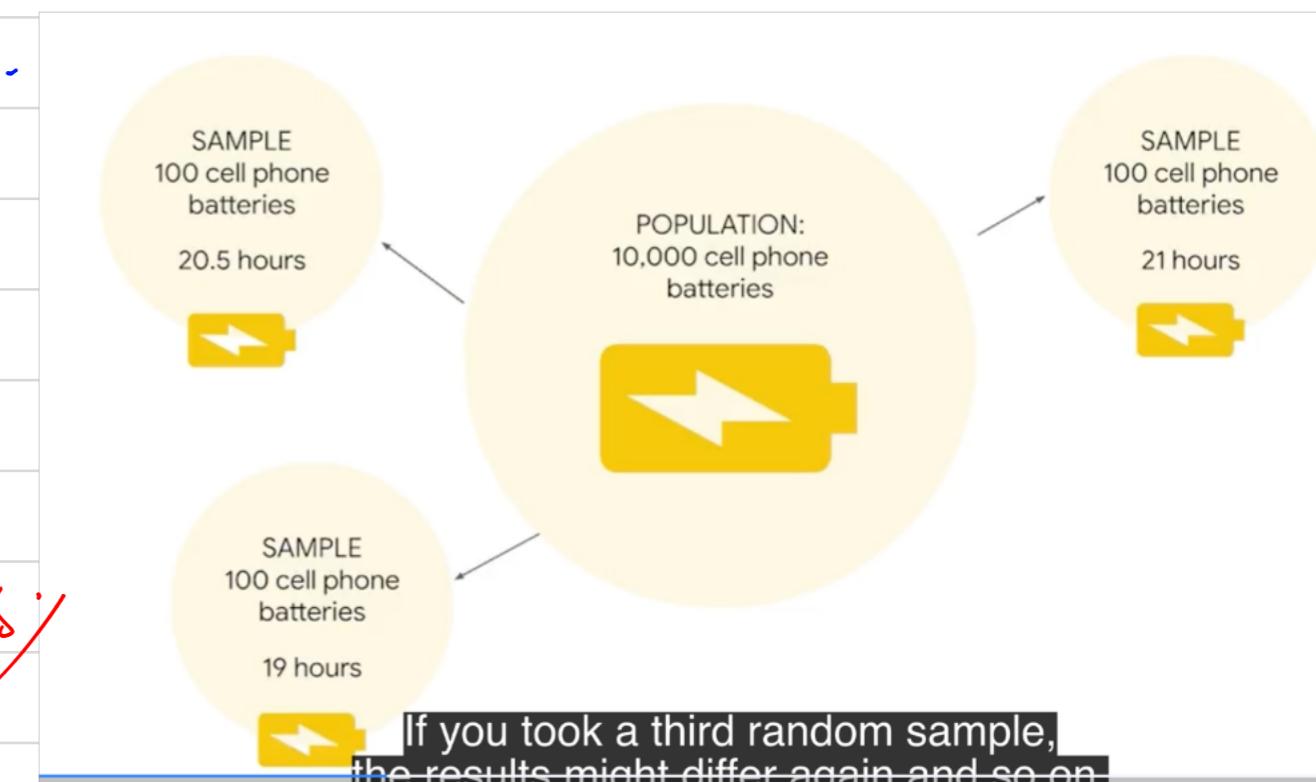
1st 100 Random Sample

2nd 100 R,

3rd

100 211

Results
differ
not results



Steps for Constructing a Confidence Interval

① Identify a sample statistic: mean

② Confidence Level: 95%

③ Margin Error = $Z\text{-Score} \times SE$

$Z\text{-Score} = \frac{\text{distance}}{\text{between}} \downarrow$
low point.

Standard Errors of the Mean:

$$SE \Rightarrow \frac{\sigma}{\sqrt{n}}$$

If σ when std of population is Known, otherwise ' s '

$$SE \Rightarrow 0.15$$

$$\text{Margin Error} = Z\text{-Score} \times SE \Rightarrow 1.96 \times 0.15 \Rightarrow 0.294$$

(4) Calculate the Interval: $[20:12, 20:48]$ → means 20 hrs 12 minutes
to 20 hrs 48 minutes.

$$20.5 + 0.294 \Rightarrow 20.794 \rightarrow \underline{\text{upper limit}}$$

$$20.5 - 0.294 \Rightarrow 20.206 \rightarrow \underline{\text{lower limit}}$$

(5) lower limit of CI = 20:12

20:12 \rightarrow 20:00
(Hypothesized).
(Cloud).

What if MEAN of Marketing? doesn't satisfy \rightarrow for that we do

Some Sample who using gg].

$[20:07, 20:53]$

Still above 20 hours

in out the
confidence interval

Confidence Intervals

* Confidence level = 95%.

$$[20:12, 20:48] = 36 \text{ min.}$$

* Confidence level = 99%.

$$[20:07, 20:53] = 46 \text{ min.}$$

What if σ is unknown? If we don't know the σ we have to use different method of calculation?

Source: The Power of Statistics Offered by Google Gurseva.