# THE CENTRAL LIMIT THEOREM

Viraj Dsouza

10 May 2024

Mathematics of Uncertainty

# IN SIMPLE WORDS...

- **Regardless** of the Population distribution, as the Sample size **increases**,
    - Sample Mean tends to Normally Distribute around the Population Mean, and
    - Sample Standard Deviation shrinks

- Understanding this empirically (By Simulating Random samples)

# POPULATION VS SAMPLE
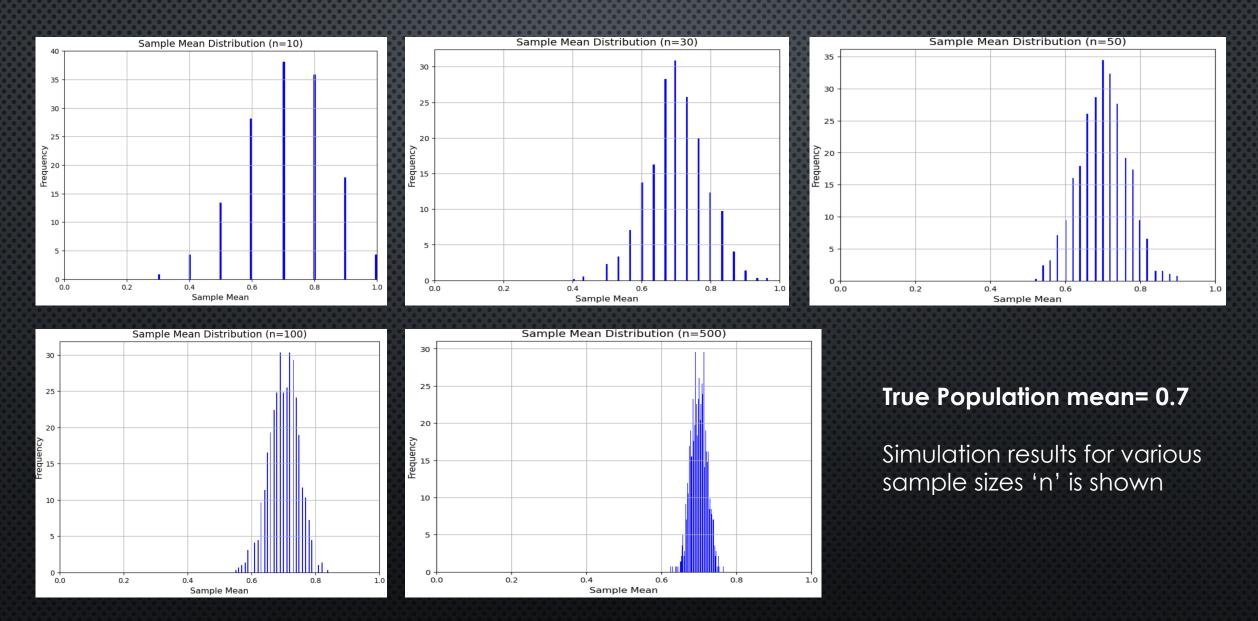
**A POPULAR ARTICLE CLAIMS:**

OUT OF THE POPULATION OF 10 MILLION PEOPLE WHO CAN VOTE, 70% SUPPORT PARTY **A** AND 30% SUPPORT PARTY **B**

GIVEN YOU HAVE ACCESS TO THIS POPULATION (BUT LIMITED RESOURCES), HOW WOULD YOU VERIFY THE ABOVE CLAIM?

RANDOM SAMPLING

**ESTIMATE** THE SUPPORT FOR PARTY A → LAW OF LARGE NUMBERS

# RANDOM SAMPLE SIMULATIONS



**True Population mean= 0.7**

Simulation results for various sample sizes 'n' is shown

# FROM THE SIMULATIONS...

**Variability** in Sampling Distribution Decreased as Sample Size Increased

Estimate from a Larger Sample size  → More Accurate **(Low sampling error)**


**Mean** of Sample Distribution Peaked close to the Population mean, for sufficiently large Sample sizes

Thoughtful Data collection → Randomizing samples (Minimizes **Bias)**


**In real world, Sampling Distributions are almost Never observed**

## THE CENTRAL LIMIT THEOREM

LET $(X_1, X_2, \dots, X_n)$ BE A RANDOM SAMPLE OF SIZE $'n'$ FROM A DISTRIBUTION WITH A FINITE MEAN $(\boldsymbol{\mu})$ AND A FINITE VARIANCE $(\boldsymbol{\sigma^2})$. IF $'n'$ IS SUFFICIENTLY LARGE, THEN $\overline{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$ IS APPROXIMATELY A

$$N\left(\boldsymbol{\mu}, \frac{\sigma^2}{n}\right) \text{ NORMAL DISTRIBUTION}$$
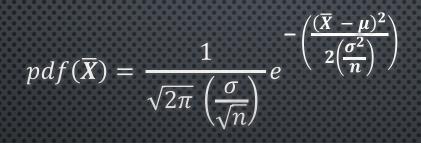
- SAMPLES MUST BE INDEPENDENT; $X_i's$ ARE INDEPENDENT RANDOM VARIABLES

- $n$ : DEPENDS ON THE POPULATION DISTRIBUTION

- MEAN OF $\overline{X}$ $(say\ \mu_{\overline{X}})\ tends\ to\ \mu\ if\ CLT\ conditions\ hold$

- SD OF $\overline{X}$ $(say\ \sigma_{\overline{X}})\ tends\ to\ \frac{\sigma}{\sqrt{n}}$ ( STANDARD ERROR) $if\ CLT\ conditions\ hold$
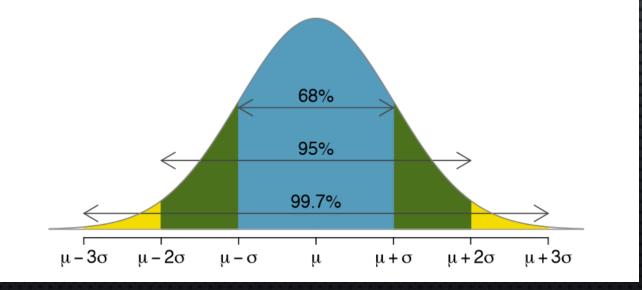
**Point Estimates**

# NORMAL DISTRIBUTION (RECALL)

If $\bar{X}$ is a $N\left(\mu, \dfrac{\sigma^2}{n}\right)$ distribution. The pdf of $\bar{X}$:

$$pdf(\bar{X}) = \frac{1}{\sqrt{2\pi}\left(\frac{\sigma}{\sqrt{n}}\right)} e^{-\left(\frac{(\bar{X}-\mu)^2}{2\left(\frac{\sigma^2}{n}\right)}\right)}$$

$N(\mu, \sigma^2)$:



68%

95%

99.7%

$\mu-3\sigma$    $\mu-2\sigma$    $\mu-\sigma$    $\mu$    $\mu+\sigma$    $\mu+2\sigma$    $\mu+3\sigma$

$$P\left(\mu - \textcolor{red}{1.96}\left(\frac{\sigma}{\sqrt{n}}\right) \leq \bar{X} \leq \mu + 1.96\left(\frac{\sigma}{\sqrt{n}}\right)\right) \approx 0.95$$

**Z-score**

# THE CLT IS VERY POWERFUL – QUITE LITERALLY

- **If CLT conditions hold, Sampling Distribution (of $\bar{X}$) Is** $N\left(\mu, \dfrac{\sigma^2}{n}\right)$

- Population Mean can Be **estimated** using sample mean ($\mu_{\bar{X}}$), but There is some $\sigma_{\bar{X}}$
  - How confident are you about the estimate based on samples You collected From the Population ?
    - **Example:** We are 95% of the times confident that the Population mean is between $(\mu - 1.96\,\sigma_{\bar{X}},\ \mu + 1.96\,\sigma_{\bar{X}})$

CLT illustrates LAW of Large Numbers (LLN)

Interval Estimate

# LONG HISTORY SHORT

- CLT IS JUST A POWERFUL EXTENSION OF LLN.

- LLN: CARDANO (16$^{TH}$ CENT.) → BERNOULLI, DE MOIVRE(18$^{TH}$ CENT.) → POISSON (19$^{TH}$ CENT.) → MARKOV, CHEBYSHEV, KOLMOGOROV, BOREL… (LATE 19$^{TH}$ AND 20$^{TH}$ CENT.)

- MONTE CARLO METHODS (1940)

CAREFUL: GAMBLER'S FALLACY

INAPPROPRIATE USE OF LLN/CLT MAY LEAD TO SERIOUS TROUBLE

# CLT IN ACTION

- Political Polling, Product-Market Fit → Public opinion surveys

- Clinical Trials

- Forecasting Weather, Stock market

- Physics: Measurement errors, Diffusion equation (Recall: Random Walk)…

# REFERENCES

- OPENINTRO STATISTICS. AVAILABLE [HERE](#)

- HOGG, R.V. , TANIS, E. AND ZIMMERMAN, D. (2015) PROBABILITY AND STATISTICAL INFERENCE. 9TH EDITION, PEARSON, UPPER SADDLE RIVER.

- FOR CLT APPLICATIONS: [INVESTOPEDIA](#)