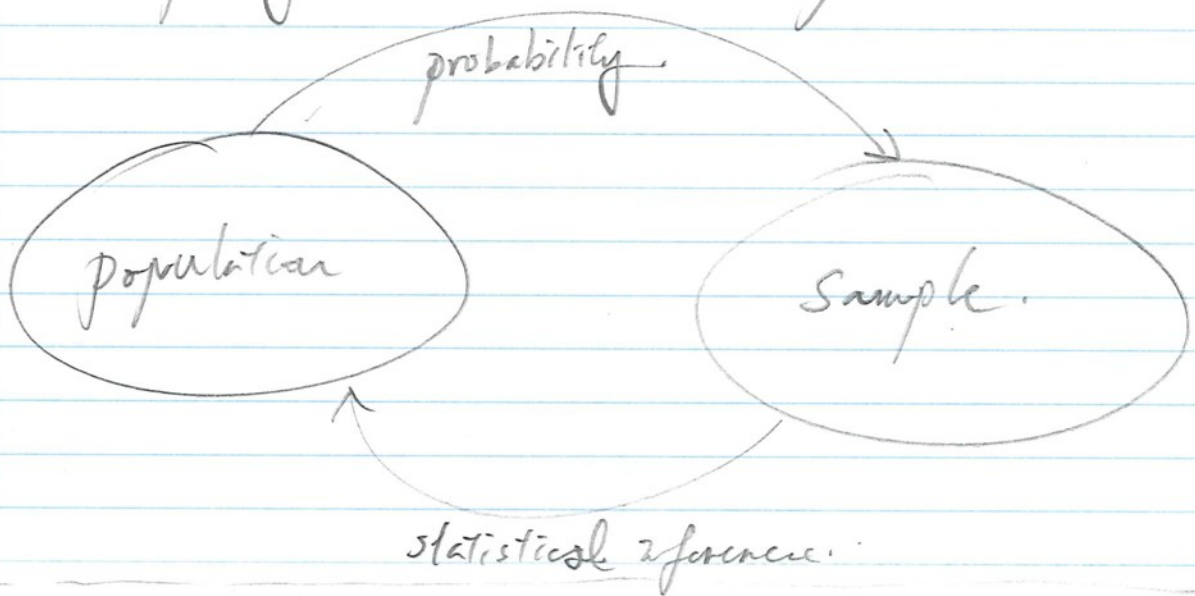


Nov 19, 2025

Sampling / Statistical inference.



★ Goal: to infer something about the population at large using a random sample.

a lot of math only works if the sample is randomly drawn from the population.

Sample version → Population version

Mean

$$\hat{\mu} = \frac{1}{n} \sum X_i = \bar{X} \quad \rightarrow \quad \mu = \mathbb{E}(X)$$

"sample mean" "population mean"

Variance

Version 1: Biased / Naïve

$$s^2 = \frac{1}{n} \sum (X - \bar{X})^2$$

$$\sigma^2 = \mathbb{E}[(X - \mathbb{E}(X))^2]$$

Version 2: Unbiased / Corrected

$$= \mathbb{E}(X^2) - \mathbb{E}(X)^2$$

$$s^2 = \hat{\sigma}^2 = \frac{1}{n-1} \sum (X - \bar{X})^2$$

$$= \text{Var}(X) / V(X).$$

var(x) R's default behavior

Covariance $\widehat{\text{Cov}}(X, Y)$

$$= \frac{1}{n-1} \sum (x - \bar{x})(y - \bar{y}).$$

↑

$R = \widehat{\text{cov}}(x, y)$

$$\begin{aligned} \text{Cov}(X, Y) &= E[(X - E(X))(Y - E(Y))] \\ &= E(XY) - E(X)E(Y) \end{aligned}$$

Statistical inference:

- We want to know
 1. A certain characteristic of the population (eg. average life expectancy)

Solution: use sample statistics to infer

→ Sample mean \rightarrow population mean
 \bar{x} \rightarrow μ

⇒ 2. How confident we are, in our estimate, given that we only have 1 sample.

Solution: Central limit theorem (CLT) + std. errors / Confidence intervals.

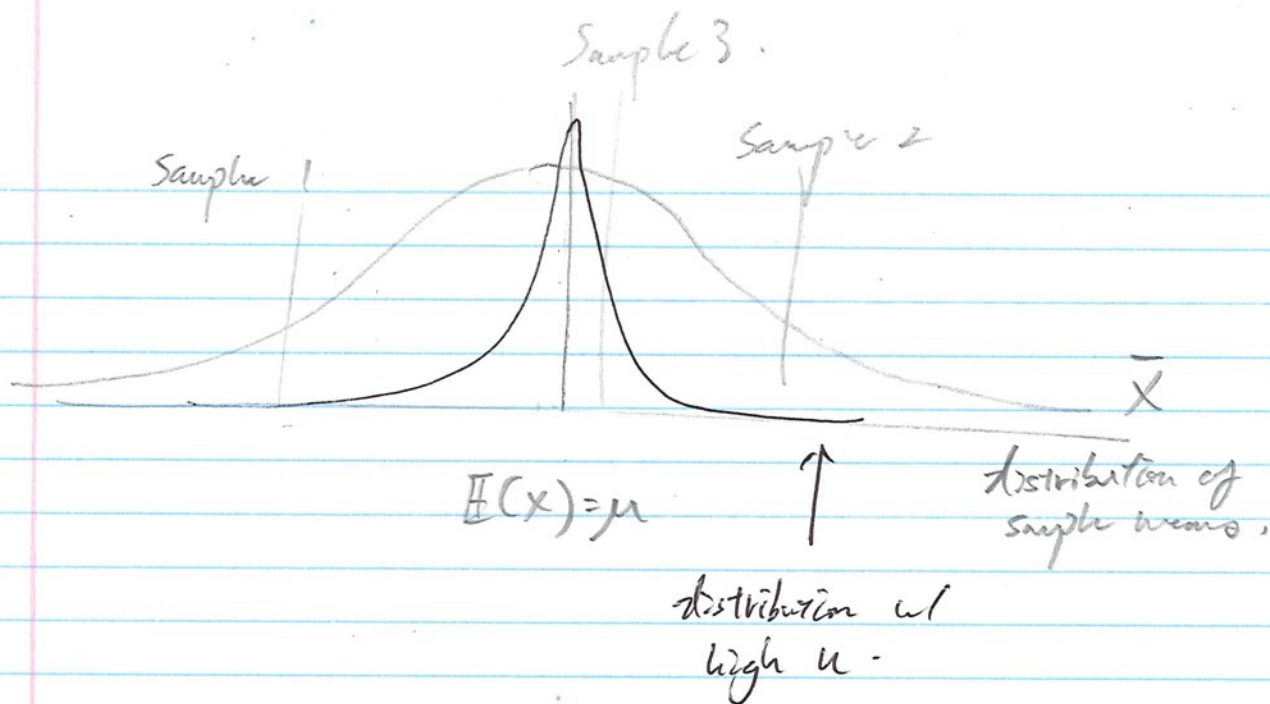
Central Limit Theorem:

Let $\mu = E(X)$ and $\sigma^2 = \text{Var}(X)$,

the sample statistics \bar{X} is distributed acc. to a normal distribution with mean μ and variance $\frac{\sigma^2}{n}$

" $\sqrt{n}(\bar{X} - \mu) \xrightarrow{d} N(0, \sigma^2)$ "

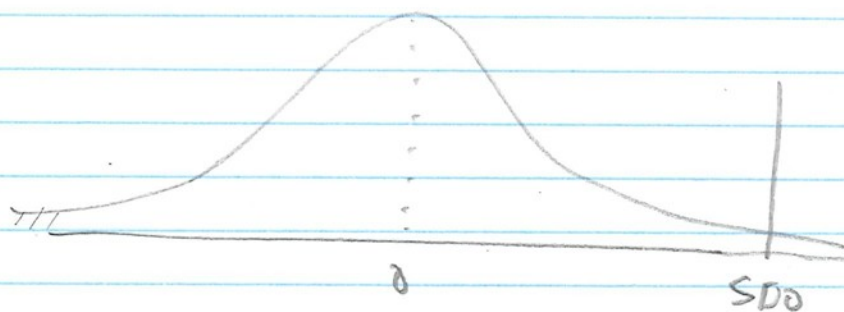
"Standard error" = std. dev. of this dist = $\sqrt{\frac{\sigma^2}{n}}$



Quantifying "confidence"

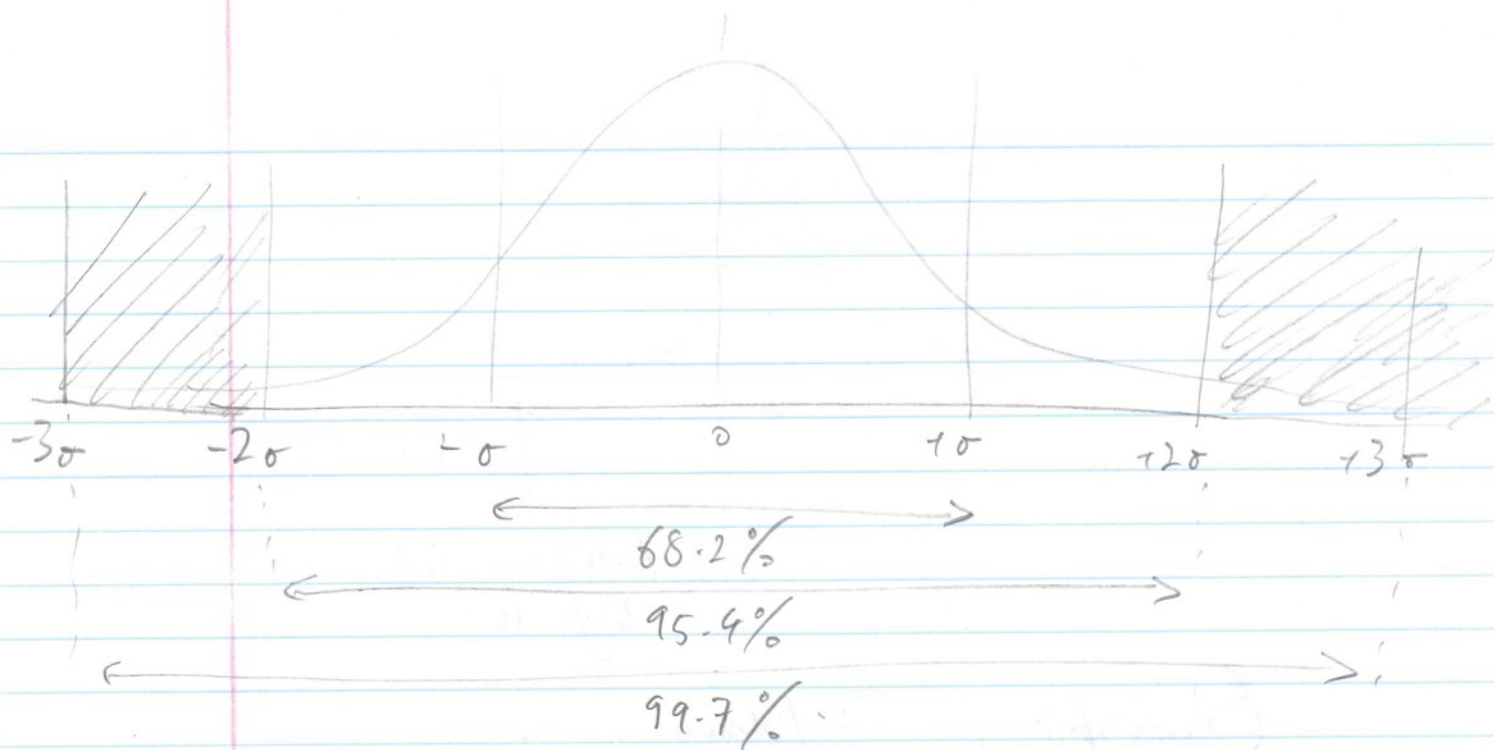
→ eg. SDD - we want to know whether the SDD we have is meaningfully different from 0.

→ Apply CLT: if the true effect is 0, then the sample SDDs should be normally distributed w/ mean 0, and variance $\frac{\sigma^2}{n}$.



- Confidence: how lucky do we have to be to observe an SDD that far in the dist. by sheer luck.

→ $= \Pr(Z > SDD)$ \leftarrow we know this because the norm dist is standard!!



Convention: we treat a statistic as "significant" if it falls outside the 95% threshold, i.e.

$$SD \geq 1.96 \text{ s.e. or } SD \leq -1.96 \text{ s.e.}$$

\hookrightarrow $\approx 5\%$ false positive rate,

Varianth 5: Confidence intervals:

$$95\% = [\text{Statistic} - 1.96 \text{ s.e.}, \text{Statistic} + 1.96 \text{ s.e.}]$$

\hookrightarrow if returned "crosses" 0, the statistic is not significant.

P-value: probability of getting a number that big by chance.

\rightarrow significant if $p < 0.05$,
(i.e. 5%).