

Lecture 14

AMS 131

Contaminated water supply:

X = arsenic concentration } both ≥ 0 difficult
 Y = lead concentration

$R = \frac{Y}{X+Y}$ proportion of contamination due to lead. $E(R) = E\left(\frac{Y}{X+Y}\right)$

Randomly sample n pairs (X_i, Y_i) from the joint PDF of (X, Y) and calculate:

$R_i = \frac{Y_i}{X_i + Y_i}$ and $\bar{R}_n = \frac{1}{n} \sum_{i=1}^n R_i$ good Monte Carlo estimate (simulation) of $E(R)$

How big does n need to be to achieve a desired accuracy?

$$|R_i| = \left| \frac{Y_i}{X_i + Y_i} \right| \leq 1 \rightarrow v(R_i) \leq \frac{1}{4}$$

CLT says that the dist. of \bar{R}_n will be close to normal for large n , with mean $E(R_i)$ and variance:

$$\frac{v(R_i)}{n} \leq \frac{1}{4n} \quad \text{PDF of } \bar{R}_n \quad \text{SD} \leq \frac{1}{2\sqrt{n}}$$

by CLT \downarrow n large

standard units $\frac{x - 0}{\frac{I}{SD}}$

suppose we want \bar{R}_n to differ from $E(R_i)$ by no more than tolerance T with probability $(1 - \alpha)$, so

$$\frac{1}{SD} \geq 2\sqrt{n} \rightarrow -\frac{I}{SD} \leq 2T\sqrt{n}$$

use inverse calculator


$$x = \Phi^{-1}\left(\frac{\alpha}{2}\right) = \frac{[E(R_i) - T] - E(R_i)}{SD} = \frac{-I}{SD} \leq 2T\sqrt{n}$$

$$n \geq \left[\frac{\Phi^{-1}\left(\frac{\alpha}{2}\right)}{2T} \right]^2 \quad \text{say } \alpha = 0.02 \quad \text{and } T = 0.005 \quad n \geq \left[\frac{-2.326}{2(0.005)} \right]^2 = 54,119$$

simulation replications

AMS 131 Lecture 14 (cont.)

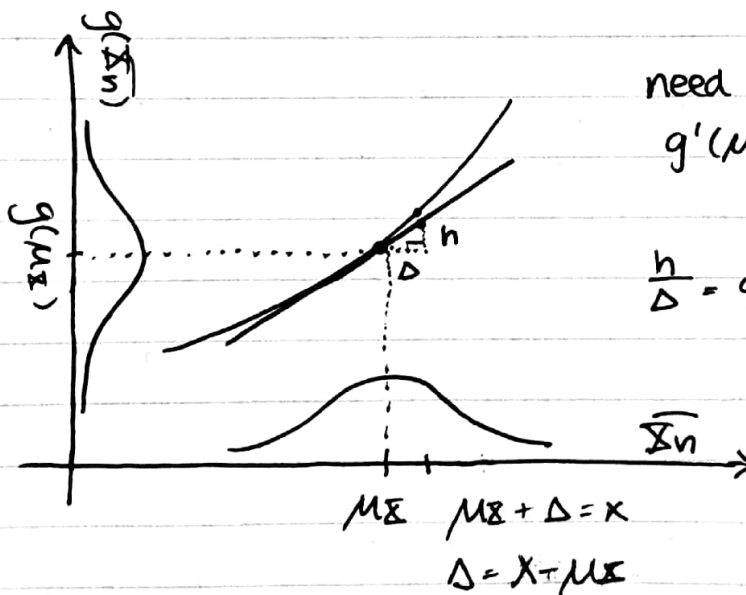
The Delta Method: the CLT says that if $X_i \sim \text{i.i.d. any dist. with finite mean and variance}$, then the dist of $\frac{\bar{X}_n - \mu_X}{\sigma_X/\sqrt{n}}$ for large n is approx. standard normal

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \quad \text{SD } \frac{\sigma_X}{\sqrt{n}} \quad \bar{X}_n \approx N(\mu_X, \frac{\sigma_X^2}{n})$$


If $g(x)$ is a sufficiently nice function, is there a comparable result for $g(\bar{X}_n)$
 \rightarrow Taylor series based approach: delta method

\bar{X}_n should be close to μ_X for large n (the weak law of large numbers)

$$X = \mu_X : g(\bar{X}_n) \approx g(\mu_X) + g'(\mu_X)(\bar{X}_n - \mu_X)$$



need to assume $g'(\mu_X) \neq 0$

$$\begin{aligned} \frac{h}{\Delta} &= g'(\mu_X) & g(x) &\approx g(\mu_X) + h \\ & & &= g(\mu_X) + g'(\mu_X) \cdot \Delta \\ & & &= g(\mu_X) + g'(\mu_X)(X - \mu_X) \end{aligned}$$

$$\begin{aligned} g(\bar{X}_n) &\approx g(\mu_X) + g'(\mu_X)(\bar{X}_n - \mu_X) \\ E[g(\bar{X}_n)] &= E[g(\mu_X) + g'(\mu_X)(\bar{X}_n - \mu_X)] \\ &= g(\mu_X) + g'(\mu_X) \underbrace{[E(\bar{X}_n) - \mu_X]}_0 \end{aligned}$$

$$E[g(\bar{X}_n)] \approx g(\mu_X) = g[E(\bar{X}_n)]$$

$$\begin{aligned} V[g(\bar{X}_n)] &\doteq V[g(\mu_X) + g'(\mu_X)(\bar{X}_n - \mu_X)] \\ &= [g'(\mu_X)]^2 \cdot V[\bar{X}_n - \mu_X] \end{aligned}$$

$$V[g(\bar{X}_n)] = [g'(\mu_X)]^2 V(\bar{X}_n) = [g'(\mu_X)]^2 \frac{\sigma_X^2}{n}$$

This works for any R.V. with $V < \infty$ as $g'(\mu_X) \neq 0$

\forall any R.V. with finite variance σ_V^2 (and \therefore

finite mean μ_V), $W = g(V) \rightarrow$

$$E(W) = g(\mu_V) \text{ and } V(W) = [g'(\mu_V)]^2 \sigma_V^2$$

provided $g'(v)$ is continuous and $g'(\mu_V) \neq 0$

if $V \sim \text{Normal}$, then $W = g(V) \sim \text{Normal}$ also

X_i = time customer i waits from reaching the front of the line until served

- Usually varies by day of week and time of day, so pick a single time slot (ex. M 10:00-10:15 AM) and observe X_i week to week in that time slot

i indexes week, so $\{X_i, i=1, 2, \dots, n\}$ forms

a stochastic process with fixed $E(X_i) = \mu_X > 0$

· stationary: non-time-varying mean $V(X_i) = \sigma_X^2$
non-time-varying variance

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \quad \text{for large } n \quad \text{rate of service: } g(\mu_X) = \frac{1}{\mu_X}, \quad g(\bar{X}_n) = \frac{1}{\bar{X}_n}$$

$$E(\bar{X}_n) = \mu_X \quad g(x) = \frac{1}{x} \quad g'(x) = -\frac{1}{x^2}$$

$$V(\bar{X}_n) = \frac{\sigma_X^2}{n} \quad g'(\mu_X) = -\frac{1}{\mu_X^2} \quad \bar{X}_n \sim \text{Normal by CLT}$$

Δ -method: $g(\bar{X}_n) = \frac{1}{\bar{X}_n} \sim \text{Normal}$ with,

says

$$g(\mu_X) = \frac{1}{\mu_X} \quad \text{variance} = \frac{\sigma_X^2}{n\mu_X^4} \quad [g'(\mu_X)]^2 = \frac{1}{\mu_X^4} \neq 0$$

Lecture 14 (cont.)

$(X_i | \lambda) \stackrel{\text{IID}}{\sim}$ Exponential(λ) reasonable for waiting times

$$E(X_i) = \mu_X = \frac{1}{\lambda} \quad V(X_i) = \sigma_X^2 = \frac{1}{\lambda^2}$$

$$\text{PDF: } f_{X_i}(x_i | \lambda) = \lambda e^{-\lambda x_i} \mathbb{I}(x_i > 0)$$

So $\frac{1}{\bar{X}_n}$ should be approx. Normal for large n with:

$$\text{mean } \frac{1}{\lambda} = \mu_X \quad \text{and SD } \frac{\sigma_X}{\mu_X^2 \sqrt{n}} = \frac{(\frac{1}{\lambda})}{(\frac{1}{\lambda})^2 \sqrt{n}} = \frac{\lambda}{\sqrt{n}}$$

$\mathbb{Y}_1, \mathbb{Y}_2, \dots$ sequence of R.V., F^* continuous CDF
 θ real number, $a_n, a_2 \rightarrow \infty$ positive sequence
 $g(\cdot)$ a real-valued function of a R.V. such that $g'(\cdot)$ is continuous and $g'(\theta) \neq 0$, then if $a_n (\mathbb{Y}_n - \theta) \xrightarrow{D} F^*$,

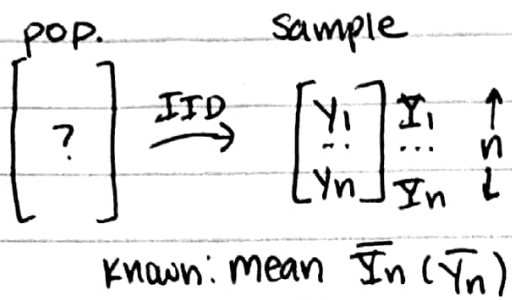
$$a_n \left[\frac{g(\mathbb{Y}_n) - g(\theta)}{|g'(\theta)|} \right] \xrightarrow{D} F^* \text{ also}$$

typical application: X_1, X_2, \dots IID

$$\bar{\mathbb{Y}}_n = \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \quad \theta = \mu_X \quad a_n = \frac{\sqrt{n}}{\sigma_X}$$

$F^* = \Phi$, the standard normal CDF

if $\frac{\bar{X}_n - \mu_X}{\sigma_X / \sqrt{n}} \rightsquigarrow N(0,1)$ then $\frac{g(\bar{X}_n) - g(\mu_X)}{|g'(\mu_X)| \sigma_X / \sqrt{n}} \rightsquigarrow N(0,1)$ also



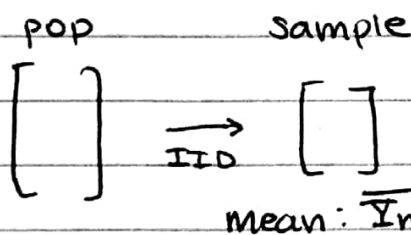
with unknown mean, SD, and PDF of the population

This is statistical inference: we want to draw inferential conclusions about μ

(Based on the data, μ is around 1. give or take 2. and I'm highly confident μ is between 3. and 4.)

confidence interval (C.I.)

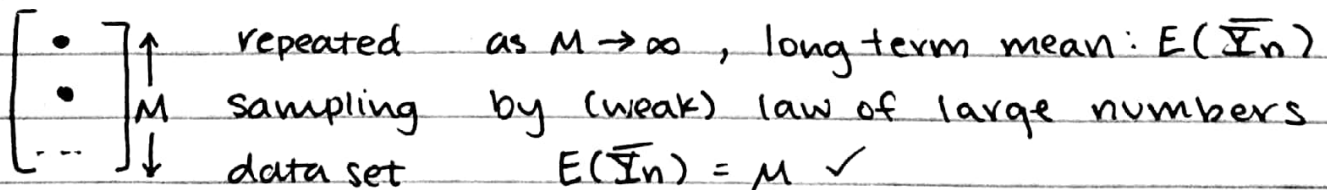
1. \bar{Y}_n (\bar{Y}_n)
 2. $\frac{\sigma}{\sqrt{n}}$
 3. $\bar{Y}_n - 2\frac{\sigma}{\sqrt{n}}$
 4. $\bar{Y}_n + 2\frac{\sigma}{\sqrt{n}}$
- (or t^{-1} number)



new hypothetical sample with a different \bar{Y}_n

How probable is it that \bar{Y}_n will differ from μ by more than _____?

possible \bar{Y}_n

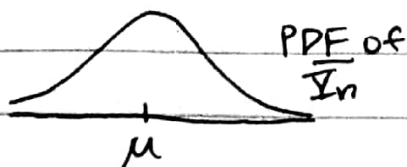


as $M \rightarrow \infty$, long term SD: $\sqrt{V(\bar{Y}_n)} \rightsquigarrow V(\bar{Y}_n) = \frac{\sigma^2}{n}$

long term PDF:

if $\hat{\theta}_n$ estimates θ , $SD(\hat{\theta}_n) \triangleq$ standard error (SE) of $\hat{\theta}_n$

$SE(\bar{Y}_n) = SD(\bar{Y}_n)$

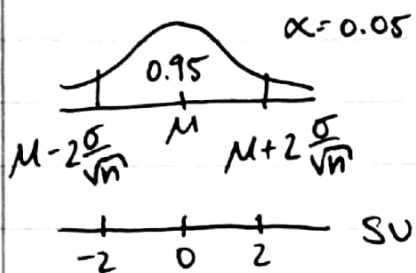


n large enough for CLT to apply

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



pick $0 < \alpha$, but small
pretend σ known (step 1)



frequentist.

$$P_F \left(\mu - 2 \frac{\sigma}{\sqrt{n}} < \bar{Y}_n < \mu + 2 \frac{\sigma}{\sqrt{n}} \right) \doteq 95\%$$

Neyman's confidence
trick

$$\mu < \bar{Y}_n + 2 \frac{\sigma}{\sqrt{n}}$$

$$\mu > \bar{Y}_n - 2 \frac{\sigma}{\sqrt{n}}$$

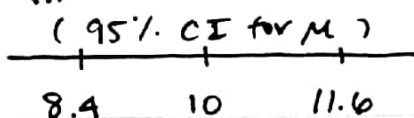
$$P_F \left(\bar{Y}_n - 2 \frac{\sigma}{\sqrt{n}} < \mu < \bar{Y}_n + 2 \frac{\sigma}{\sqrt{n}} \right) \doteq 95\%$$

\therefore lets define $(\bar{Y}_n \pm 2 \frac{\sigma}{\sqrt{n}})$ to be a
95% confidence interval for μ

$$\bar{Y}_n = 10$$

$$n = 25 \quad \bar{Y}_n \pm 2 \frac{\sigma}{\sqrt{n}} = 10 \pm 2 \left(\frac{4}{5} \right) = (8.4, 11.6)$$

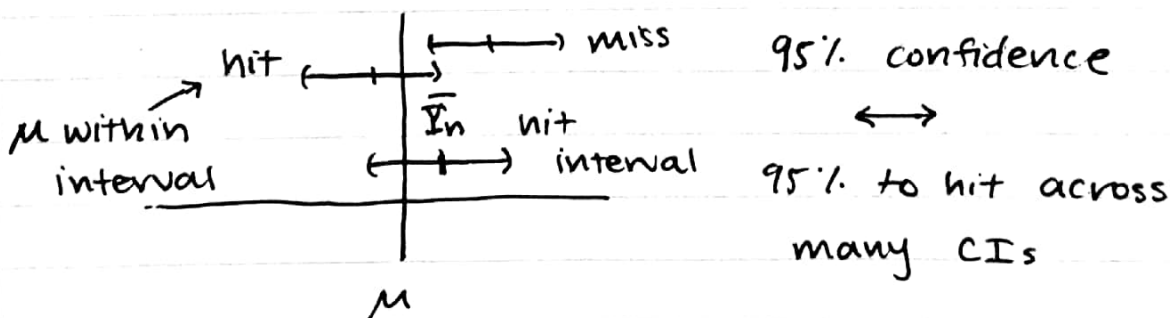
$$\sigma = 4$$



easy to conclude that $P_F(8.4 < \mu < 11.6) \doteq 95\%$.

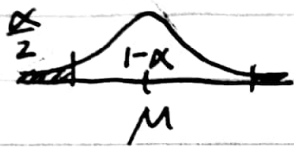
but this is wrong, μ is a fixed unknown constant

\therefore no frequentist probability statements can
be made about it: $P_F(8.4 < \mu < 11.6) = \text{undef.}$



We will not know if our 95% CI is a
hit or miss, only way to know is if
you know μ - if you knew μ , you wouldn't
need a CI in the first place.

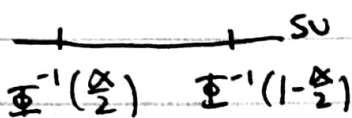
our confidence is in the process of building the CI, not
the outcome (the CI itself)



100(1- α)% CI:

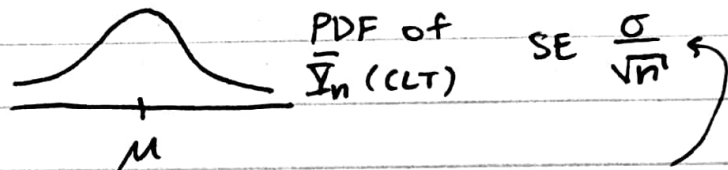
$$\bar{Y}_n \pm z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$$

100(1- α)% CI have a 100 $\cdot\alpha$ % false-positive (miss) rate




95% CI: too low for good, replicable science
 (0.05) false positive rate \rightarrow (0.005), 99.5% CI

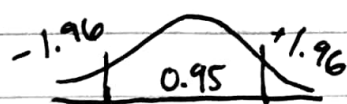
σ unknown (step 2)



$\bar{Y}_n \pm z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$ \leftarrow now use S_n , the SD of the sample, not normal anymore
 long term SD $\sqrt{\hat{V}(\bar{Y}_n)} = \frac{S_n}{\sqrt{n}}$

n large $\rightarrow z_{1-\frac{\alpha}{2}}$ approx. OK

n small  $\frac{\bar{Y}_n - \mu}{S_n/\sqrt{n}}$ CLT



\downarrow t curve

t_{n-1} degrees of freedom

$n \uparrow$, PDF \rightarrow normal, $t_n \rightarrow \Phi$ standard normal

Lecture 14 (cont.)

T-S case study

\bar{X} = # of T-S babies in family of 5, both parents are carriers

$P(\text{T-S baby}) = \frac{1}{4} = p$ $\bar{X} \sim \text{Binomial}(n, p)$

$T_i = \begin{cases} 1 & \text{if T-S baby} \\ 0 & \text{else} \end{cases} \quad i=1, \dots, n=5$

then $T_i \stackrel{\text{i.i.d.}}{\sim} \text{Bernoulli}(p) \quad i=1, \dots, n$ and

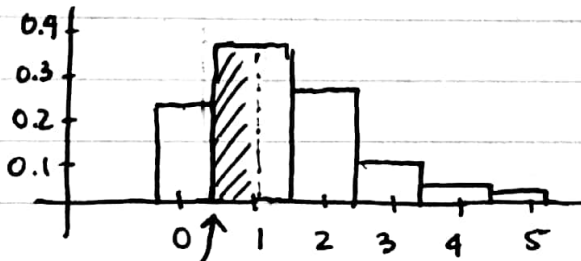
$\bar{X} = \sum_{i=1}^n T_i$ by CLT, \bar{X} should be approx. normal

$\mu_{\bar{X}} = E(\bar{X}) = np = 1.25$

$\sigma_{\bar{X}} = \sqrt{V(\bar{X})} = \sqrt{np(1-p)} \approx 0.98$

$P(\text{1 or more T-S}) = P(\bar{X} \geq 1)$

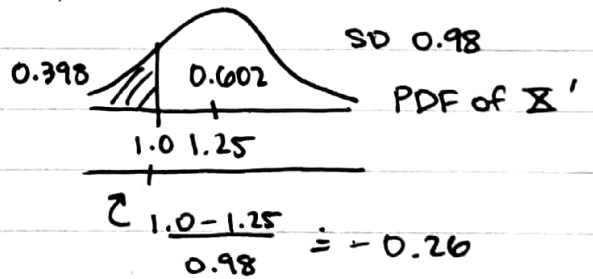
$1 - P(\text{no T-S}) = 1 - (1-p)^n \approx 0.76 = 1 - P(\bar{X} = 0)$



better approx.

naive approx.

naive normal approx. from CLT



$P(\bar{X} \geq 1) \approx 1 - P(X' < 1)$

$= 1 - 0.398$

≈ 0.602 bad approx.

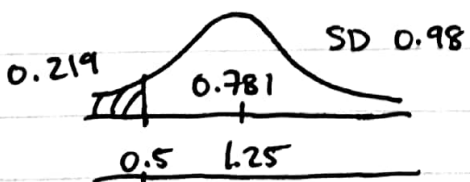
(vs. 0.76) ↗

not good approx.

because $n=5$

Continuity correction

now look at the edges of the histogram bars



$P(\bar{X} \geq 1) = 1 - P(\bar{X}' < 0.5)$

$= 1 - 0.219$

≈ 0.781

MUCH closer to 0.76

Markov Chains

A sequence of R.V.s X_1, X_2, \dots is called a stochastic process with discrete time parameter $t=1, 2, \dots$. X_1 is the initial state of the process. $X_n, n \geq 1$ is the state of the process at $t=n$.

The simplest possible discrete time stochastic process is an IID sequence of R.V.s (X_1, X_2, \dots) .

Suppose that there's a parameter θ such that $(X_i | \theta) \sim \text{IID}$ from some dist. depending on θ .

↳ machine with θ dial produces IID Bernoulli(θ)

- the process (X_1, X_2, \dots) does have memory for you if θ is unknown: the information of previous trials helps you predict future ones
- the process $(X_i | \theta)$ has no memory if θ known

An IID process $(X_i | \theta)$ is called a white-noise (stochastic) process or white noise time series.

What's the next level of complexity?

Allow X_{n+1} depend on X_n , but not X_{n-1}, X_{n-2}, \dots

the process has short term memory, 1 time period

A discrete-time stochastic process is a (first order) Markov chain if for $n=1, 2, \dots$; b , any real number; and for all possible sequences of states x_1, x_2, \dots

$$P(X_{n+1} \leq b | X_1 = x_1, \dots, X_n = x_n) \\ = P(X_{n+1} \leq b | X_n = x_n)$$

The only thing you need to simulate where the Markov chain is going next is where it is now.