

Supplemental Notes

Topics

① Central Limit Theorem (CLT)

$$- Z_n \xrightarrow{d} Z \sim N(0,1) \quad \text{random sample}$$

$$- (1-\alpha) \% \text{ C.I. for } \mu_X: \bar{X}_n \pm z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

$$Z_n = \text{STD}(\bar{X}_n) = \text{STD}\left(\sum_{k=1}^n X_k\right)$$

CLT Binomial Approximation

$$P(X=k) = \binom{n}{k} p^k (1-p)^{n-k}$$

$$\therefore X = \sum_{k=1}^n X_k \text{ independent Bernoulli } X_k$$

$$f_{X_k}(x) = p^x (1-p)^{1-x} \text{ for } x=0 \text{ or } 1$$

$$X = \sum_{k=1}^n X_k \xrightarrow{\text{CLT } d}$$

$$\text{For large } n: b(n, k, p) \approx N\left(\frac{\mu_X}{np}, \frac{\sigma_X^2}{n \cdot pq}\right) \quad \text{if } q = 1-p.$$

Rule of thumb:

$$\text{and } \begin{cases} (1) & np \geq 5 \\ (2) & nq \geq 5 \end{cases}$$

Ex: $X \sim b(10, \frac{1}{2})$

$\therefore n = 10 \quad p = \frac{1}{2}$

$\therefore np = n(1-p) = 10 \cdot \frac{1}{2} = 5.$

\therefore OK to use CLT.
(but NOT Poisson Law)

$\therefore P(3 \leq X < 6) \approx P(2.5 \leq X \leq 5.5)$
continuity correction

standardize
 $= P\left(\frac{2.5 - np}{\sqrt{npq}} \leq \frac{X - np}{\sqrt{npq}} \leq \frac{5.5 - np}{\sqrt{npq}}\right)$

$= P\left(\frac{2.5 - 5}{\sqrt{10/4}} \leq \frac{X - 5}{\sqrt{10/4}} \leq \frac{5.5 - 5}{\sqrt{10/4}}\right)$

CLT
 $\approx P(-1.58 \leq Z \leq 0.3162)$

$= \Phi(0.316) - \Phi(-1.58)$
for $Z \sim N(0, 1)$

$\approx 0.5684.$

Exact: $P(3 \leq X < 6) = 0.5683$ if $X \sim b(10, \frac{1}{2})$

$\therefore |0.5684 - 0.5683| = 0.0001$

Continuity correction:

for discrete, $P(X \leq x) = P(X < x + 1)$

for continuous, $P(X \leq x) = P(X < x)$

\therefore discrete \rightarrow continuous, evaluate midpoint $\approx P(X \leq x + \frac{1}{2})$

Note: $P(3 \leq X \leq 6) \approx P(2.5 \leq X \leq 5.5)$
continuity correction.

$P(3 \leq X \leq 6) \approx P(2.5 \leq X \leq 6.5)$

\geq, \leq wider

$>, <$ shrink

Ex: Airline overbooking

$P(\text{no show}) = 0.04$ 126 seats but 130 bookings

Q: What is probability that more than 126 passengers show up?

(i.e. 3 or fewer no-shows)

X : # no show What is $P[X \leq 3]$ i.e. $X = 0, 1, 2, \text{ or } 3$

$n = 130$ and $p = 0.04$
large sample requirement

$\therefore np = 5.2$
 $npq = 124.8$
 \therefore CLT approximation OK.

$$\therefore \sqrt{npq} = \sqrt{(130)(0.04)(0.96)} = 2.2343$$

$$\begin{aligned} P[X \leq 3] &\stackrel{\text{cont. corr.}}{\leq} P[X \leq 3.5] = P\left[\frac{X - np}{\sqrt{npq}} \leq \frac{3.5 - np}{\sqrt{npq}}\right] \\ &= P\left[\frac{X - 5.2}{2.2343} \leq \frac{3.5 - 5.2}{2.2343}\right] \\ &= P[Z \leq -0.76] \\ &\quad \downarrow \text{z-table} \\ &= 0.2236 \quad (\approx 22\%) \end{aligned}$$

Compare

Exact: $p = \sum_{k=0}^3 \binom{130}{k} (0.04)^k (0.96)^{130-k} = 0.2323 \approx 23\%$
(close)

Poisson approximation: $p = 0.04 \ll 1$, $n > 100$ $\therefore \lambda = np = 5.2$

$$P(X \leq 3) \approx \sum_{k=0}^3 \frac{e^{-5.2} (5.2)^k}{k!} = 0.23 \quad (\approx 23\%)$$

Confidence Intervals

2-sided for parameter Θ :

$$\hat{\Theta}_n \pm \text{M.O.E.}$$

M.O.E. = Margin of Error

usually involves $\sigma_{\hat{\Theta}_n} = \sqrt{V[\hat{\Theta}_n]}$
"standard error"

or its estimate $S_{\hat{\Theta}_n}$

$\hat{\Theta}_n$ usually involves \bar{X}_n, S_n^2 , or S_{xy} less confident
↓
bigger MDE

$\alpha \in [0, 1]$ (random sets $X: \Omega \rightarrow \mathcal{B}$).

Ex: $(1-\alpha)\%$ C.I. for μ_x if σ_x^2 known and

random sample and Large $n \geq 30$ CLT

$$\text{C.I.} = \bar{X}_n \pm \text{M.O.E.}$$

$$= \bar{X}_n \pm z_{\alpha/2} \left(\frac{\sigma_x}{\sqrt{n}} \right) \leftarrow \sqrt{V[\bar{X}_n]} \text{ standard error}$$

unknown σ_x^2 :

$$\text{C.I.} = \bar{X}_n \pm t_{\alpha/2}(n-1) \frac{S_n}{\sqrt{n}} \leftarrow S_n \text{ estimates } \sigma_x$$

\uparrow wider C.I.

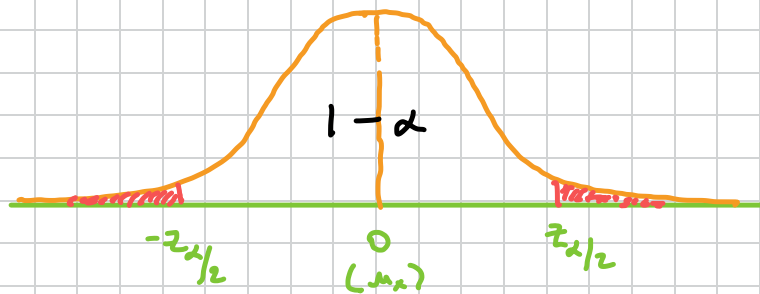
Thm: $(1-\alpha)\%$ C.I. for μ_x if σ_x^2 known and large n :

$$\bar{X}_n \pm z_{\alpha/2} \frac{\sigma_x}{\sqrt{n}}$$

memorize

(r.s.) \rightarrow IID $\rightarrow \sigma_x^2 < \infty$

Prf:



$$Z_n \stackrel{\Delta}{=} \frac{X_n - \mu_x}{\sigma_x / \sqrt{n}}$$

since r.s.

$$\therefore \text{large } n: Z_n \stackrel{\text{CLT}}{\approx} Z \sim N(0, 1)$$

$$\therefore \bar{X}_n \stackrel{\text{CLT}}{\approx} N\left(\mu_x, \frac{\sigma_x^2}{n}\right)$$

constant, i.e. $\frac{\alpha}{2}$: prob $\geq z_{\alpha/2}$

(symmetric interval)

$$\therefore 1-\alpha = \text{P.T.C.} \left[P\left(-z_{\alpha/2} \leq Z \leq z_{\alpha/2}\right) \right]$$

$$\stackrel{\text{CLT}}{\approx} P\left(-z_{\alpha/2} \leq \frac{\bar{X}_n - \mu_x}{\sigma_x / \sqrt{n}} \leq z_{\alpha/2}\right)$$

since n "large"

$$= P\left(\left\{ \omega \in \Omega: -z_{\alpha/2} \leq \frac{\bar{X}_n(\omega) - \mu_x}{\sigma_x / \sqrt{n}} \leq z_{\alpha/2} \right\}\right)$$

$$= P\left(-z_{\alpha/2} \frac{\sigma_x}{\sqrt{n}} \leq \bar{X}_n - \mu_x \leq z_{\alpha/2} \frac{\sigma_x}{\sqrt{n}}\right)$$

$$= P\left(-\bar{X}_n - z_{\alpha/2} \frac{\sigma_x}{\sqrt{n}} \leq -\mu_x \leq -\bar{X}_n + z_{\alpha/2} \frac{\sigma_x}{\sqrt{n}}\right)$$

$$= P\left(\bar{X}_n + z_{\alpha/2} \frac{\sigma_x}{\sqrt{n}} \geq \mu_x \geq \bar{X}_n - z_{\alpha/2} \frac{\sigma_x}{\sqrt{n}}\right)$$

$$= P\left(\bar{X}_n - z_{\alpha/2} \frac{\sigma_x}{\sqrt{n}} \leq \mu_x \leq \bar{X}_n + z_{\alpha/2} \frac{\sigma_x}{\sqrt{n}}\right)$$

$\therefore \bar{X}_n \pm z_{\alpha/2} \frac{\sigma_x}{\sqrt{n}}$ is the $(1-\alpha)\%$ CI for μ_x

QED.

Picking Sample Size n

Want: interval no wider than $\bar{X}_n \pm \varepsilon$ in above case

$$\varepsilon = \text{Maximal error tolerance} = \text{M.O.E.} = z_{\alpha/2} \frac{\sigma_x}{\sqrt{n}}$$

$$\therefore \varepsilon^2 = z_{\alpha/2}^2 \cdot \frac{\sigma_x^2}{n}$$

$$\therefore n = \frac{z_{\alpha/2}^2 \sigma_x^2}{\varepsilon^2}$$

MEMORIZE

Proportions: Estimate with CLT:

Q: estimate P [vote candidate]

$$X \sim b(n, p)$$

← success/failure structure

proportion $\hat{p} = \frac{X}{n}$ ← # success
← # trials
estimate for P [success]

note: $\hat{q} = 1 - \hat{p}$

(MLE of p since $X \sim b(n, p)$)

$$\textcircled{1} \quad \underline{E[\hat{p}]} = \frac{E[X]}{n} \stackrel{X \sim b}{=} \frac{n \cdot p}{n} = p \quad \forall n \quad \therefore \text{Unbiased}$$

$$\textcircled{2} \quad \underline{V[\hat{p}]} = \frac{V[X]}{n^2} = \frac{n \cdot p \cdot q}{n^2} = \frac{p \cdot q}{n} \quad \downarrow 0 \text{ as } n \uparrow \infty$$

standard error

$$\sigma_{\hat{\theta}_n} = \frac{\sqrt{pq}}{\sqrt{n}}$$

$$\text{MSE} = V + B^2$$

$$\therefore \hat{p}_n \xrightarrow{P} p$$

by m.p.d.

$\therefore \hat{p}_n$ is consistent for p .

\therefore For large n : CLT $\hat{p}_n \stackrel{d}{\approx} N\left(p, \frac{p(1-p)}{n}\right)$
(like sample-mean)

$\therefore (1-\alpha)\%$ C.I. for p :

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

Ex: Polls. ^(report) "±3% M.O.E." in the news poll

estimate $\rightarrow P[\text{vote candidate}]$

implied $\epsilon = 0.03$ at 95% C.I. ($\alpha = 0.05$)

problem: width depends on estimate \hat{p}
so take "worst case" (widest interval)

Note: $p(1-p) \leq \frac{1}{4}$ ^{since $p \in (0,1)$} $\therefore p = \frac{1}{2}$ maximizes $f(p) = p - p^2$

$$\therefore V[\hat{p}] = \frac{p \cdot q}{n} \leq \frac{1}{4n}$$

\therefore Conservative estimate $\cdot \sigma_{\hat{p}} \leq \frac{1}{2\sqrt{n}}$

$$\therefore \epsilon = \text{M.O.E.} = \frac{Z_{\alpha/2}}{2\sqrt{n}}$$

↓ solve for n

$$\therefore n = \frac{Z_{\alpha/2}^2}{4\epsilon^2}$$

pollster's sample size

\therefore For 95% confidence

Need $n \geq 1068$ samples

"
 $1-\alpha, \therefore \alpha = 0.05$
 $\frac{\alpha}{2} = 0.025$

For 99% confidence

$n \geq 1849$ samples

$\therefore \alpha = 0.005$