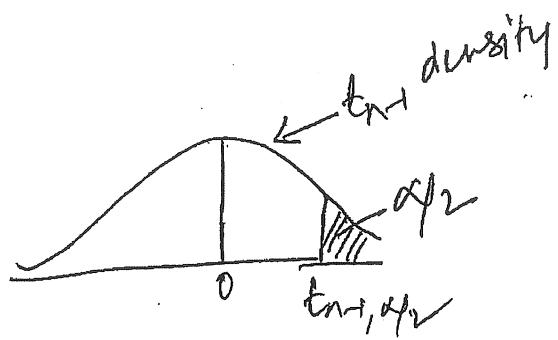
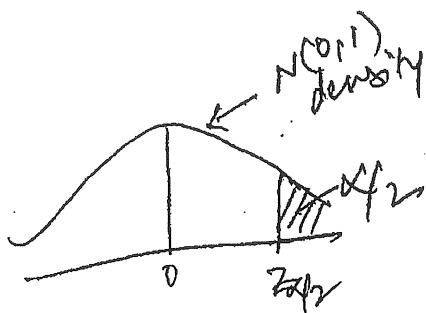
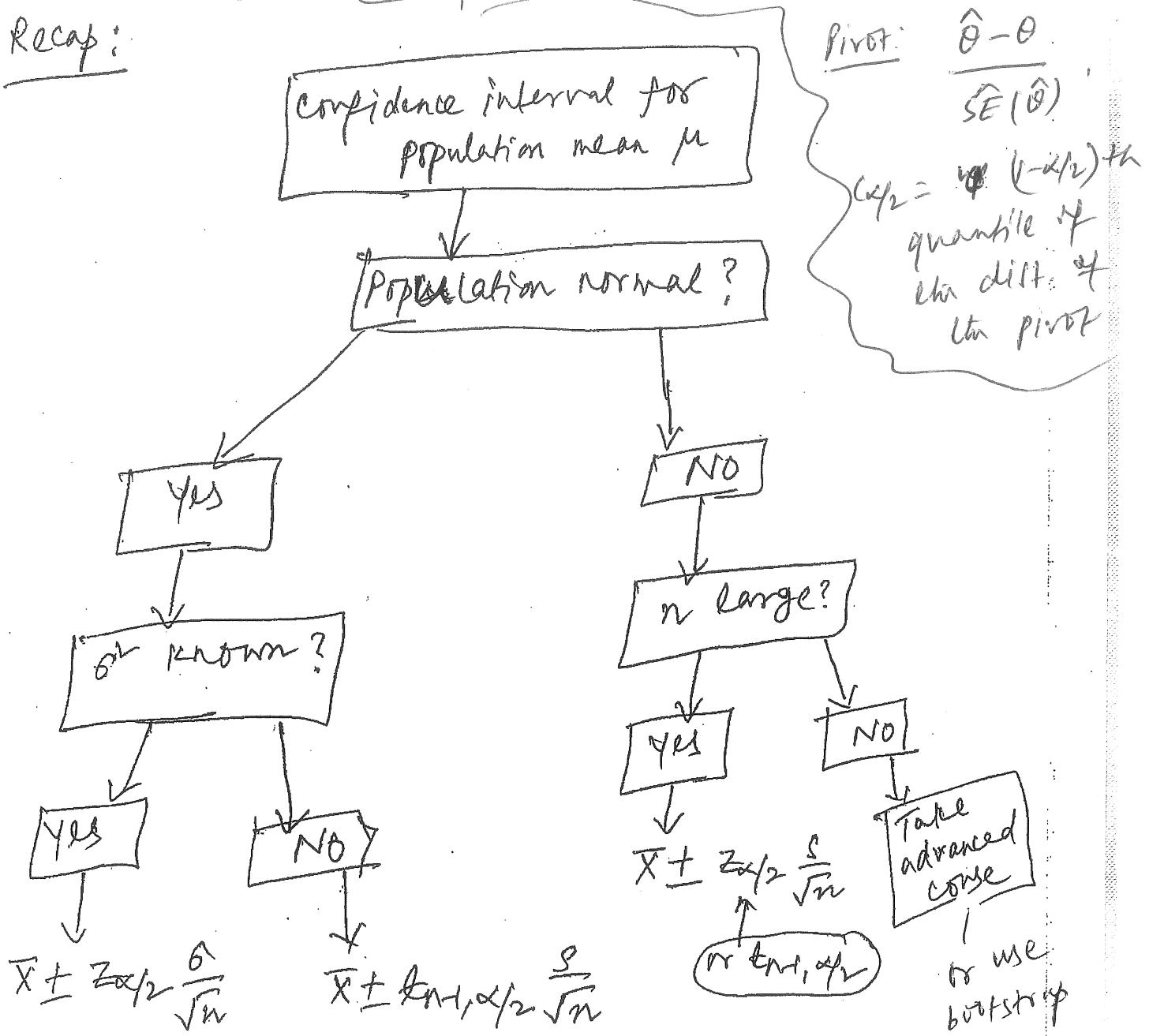


Recall: Generally: $A \cdot 100(1-\alpha)\% \text{ CI for } \theta: \hat{\theta} \pm c_{\alpha/2} \hat{SE}(\hat{\theta})$,

Recap:



Ex: If an unauthorized person accesses a computer account with the correct username and password (stolen or cracked), can this intrusion be detected? One way to do this is to compare mean time between keystrokes of the user trying to log in with that of the account owner. The intrusion is detected if there is a noticeable difference. The following times between keystrokes (in seconds) were recorded when a user typed the username and password:

$$X_1 \rightarrow X_2 \rightarrow X_3 \rightarrow X_4 \rightarrow X_5 \rightarrow X_6 \rightarrow X_7 \rightarrow X_8 \rightarrow X_9 \rightarrow X_{10} \rightarrow X_{11} \rightarrow X_{12}$$

$$0.46, 0.38, 0.31, 0.24, 0.20, 0.31, 0.34, 0.42, 0.09, 0.18, 0.46, 0.21$$

Find a 95% CI for mean time between keystrokes for the user trying to log in. Assume a normal distribution for the times.

$$X = \text{(typical) time b/w two keystrokes for the person trying to log in.}$$

σ is unknown and needs to be estimated with data.

Note:

$$95\% \text{ CI for } \mu$$

$$X \sim N[\mu, \sigma^2]$$

Goal:

$$\text{t-interval: } \bar{x} \pm t_{1-\alpha/2} \frac{s}{\sqrt{n}}$$

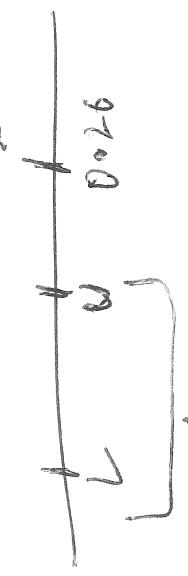
$$= \boxed{\bar{x} \pm t_{1-\alpha/2} \frac{s}{\sqrt{n}}}$$

```
x <- c(0.46, 0.38, 0.31, 0.24, 0.20, 0.31, 0.34,  
0.42, 0.09, 0.18, 0.46, 0.21)  
#> mean(x)  
#[1] 0.3  
#> sd(x)  
#[1] 0.1183216  
#>qt(0.975, 11)  
#[1] 2.200985
```

Suppose: Mean for ACE owner is 0.26

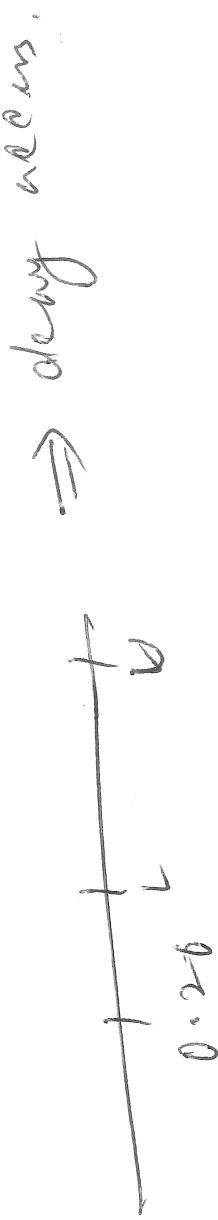
Since 0.26 is a plausible value (because
for the person who is trying to login,
access should be granted.

← ACE owner



Suppose:

Person who is trying to log in is typing
faster than the ACE owner \Rightarrow deny access



Similarly:

Large sample CI for mean μ

Recall: When n is large, an approximate $100(1 - \alpha)\%$ CI for mean μ of any population is $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N[0, 1]$

Recall: $\bar{X} \sim N[\mu, \frac{\sigma^2}{n}]$

$$\Rightarrow \bar{X} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$
$$\Rightarrow \bar{X} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

Ex: We wish to estimate the mean execution time of a program. The program was run 35 times on randomly selected inputs, and the sample mean and the sample standard deviation of the execution times were evaluated as 230 ms and 14 ms, respectively. Find a 95% CI for the true mean execution time μ .

$X = \text{(typical) execution time (ms)}$

Know: $n = 35$, use large sample CI: $\bar{X} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$

95% CI for μ :

$$= \left[225.4, 234.6 \right] \text{ ms}$$

Plausible values for μ :

$100(1-\alpha)$ % of the θ : $[L, U]$ is such that -

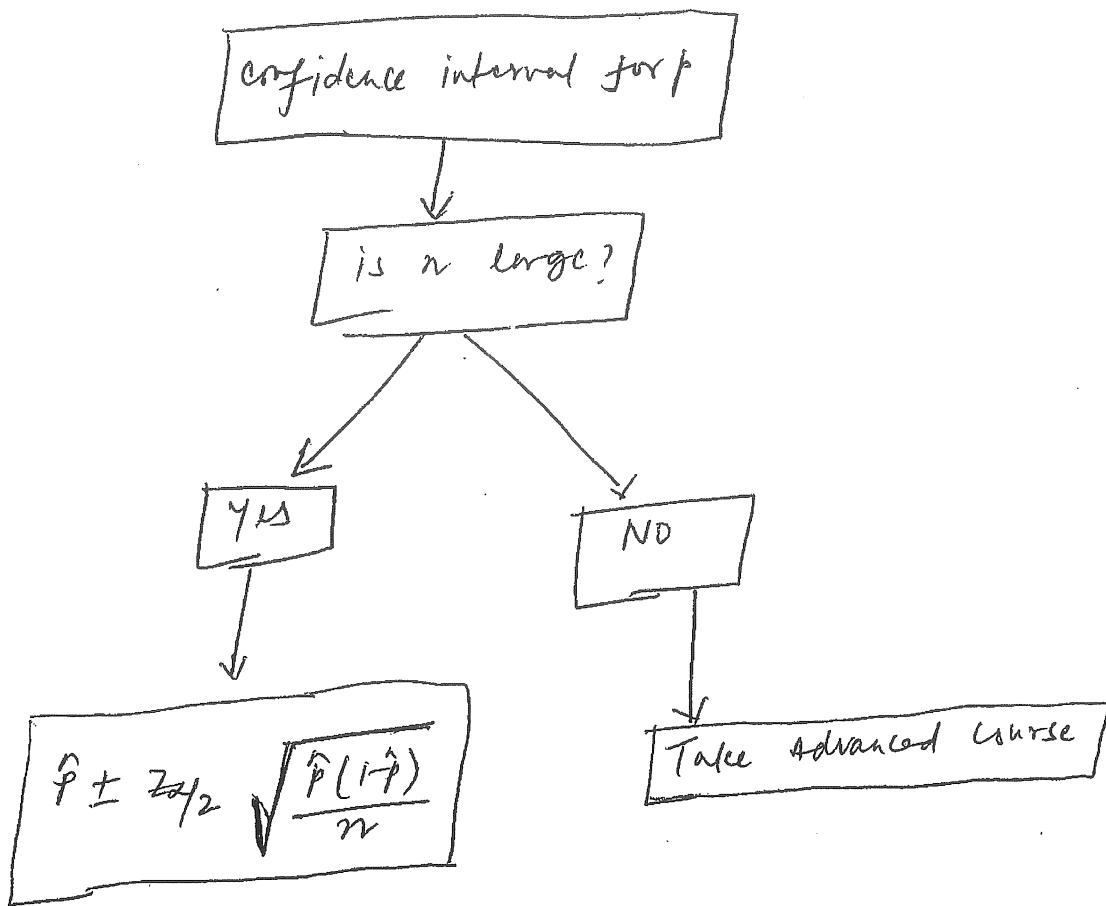
$$P\left[\frac{L \leq \theta \leq U}{U} = 1-\alpha\right] = 1-\alpha$$

- "exact"

Approx. if:

$$P\left[\frac{L \leq \theta \leq U}{U} \approx 1-\alpha\right] \approx 1-\alpha$$

$\left[\text{As } n \rightarrow \infty, P\left[L \leq \theta \leq U\right] \rightarrow 1-\alpha\right]$



Large sample CI for success proportion \hat{p}

Population: $X \sim \text{Bernoulli}(p)$, where $p = \text{proportion of successes in population}$; $p = E(X) = P[X=1]$.

Sample data: X_1, \dots, X_n . (Note: they are 0s and 1s).

Recall: Estimator for $p = \hat{p} = \frac{\bar{x}}{n}$ = proportion of successes in the sample.

Also: Estimated $\text{var}(X) = \text{estimate of } p(1-p) = \boxed{\hat{p}(1-\hat{p}) / n}$

why is
 n is
large.

Result: An approximate CI for p : $\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$

Proof:

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

Formula:

$$\Rightarrow \text{Var}[\hat{p}] = \frac{\hat{p}(1-\hat{p})}{n} \Rightarrow \text{SE}[\hat{p}] = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

Ex: From a large population of RAM chips, a random sample of 50 is taken and a test carried out on each to see whether they perform correctly. In the test, only 20 chips are found to perform correctly. Find a 95% CI for p , the true proportion of chips that perform correctly.

$$X = \begin{cases} 1, & \text{if a chip performs correctly} \\ 0, & \text{otherwise} \end{cases}$$

$X \sim \text{Bernoulli}(p)$

Now:

$$\hat{p} = \frac{20}{50} = 0.4$$

95% CI for p :

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

$\hat{p} = 0.26, 0.54$

Plausible values for p -

Choosing the sample size n

- Width of CI = $V - L = 2 z_{\alpha/2} \sqrt{\hat{p}(1-\hat{p}) / n}$
- Let w = desired CI width for $1 - \alpha$ confidence.
- Margin of error = $w/2$

- Set CI width = desired width and solve for n to get

$$\text{Solve for } n: w = \frac{2 z_{\alpha/2} \sqrt{\hat{p}(1-\hat{p})}}{n} \Rightarrow n = \left[\frac{2 z_{\alpha/2} \sqrt{\hat{p}(1-\hat{p})}}{w} \right]^2$$

- This formula involves \hat{p} , which is not known before the experiment.

- One alternative: take $\hat{p} = 0.5$ because $\hat{p}(1 - \hat{p})$ is maximum when $\hat{p} = 0.5$. This strategy will yield a conservative values of n . (The sample size will be larger than necessary.)

$$\text{Replacing } \hat{p} \text{ by } \frac{1}{2} \text{ gives: } n \approx \left[\frac{2 z_{\alpha/2}}{w} \right]^2$$

more than 50%
what is really needed.

Ex: Suppose we are planning a survey to estimate the proportion of American who approve of President Trump's job. We would like our estimate to be within 3% of the true proportion with 95% confidence. How much sample size should we take?

$$X = \begin{cases} 1, & \text{if person approves the job} \\ 0, & \text{otherwise} \end{cases}$$

$1 - \alpha = 0.95$

Want: $| \hat{p} - p | \leq 0.03$

↓
desired margin of error

$$n = \frac{2(0.03)^2}{0.06} = 0.06$$

$\Rightarrow n \approx 1068.$

$$\left[\frac{1.96}{0.06} \right]^2 \approx 1068.$$

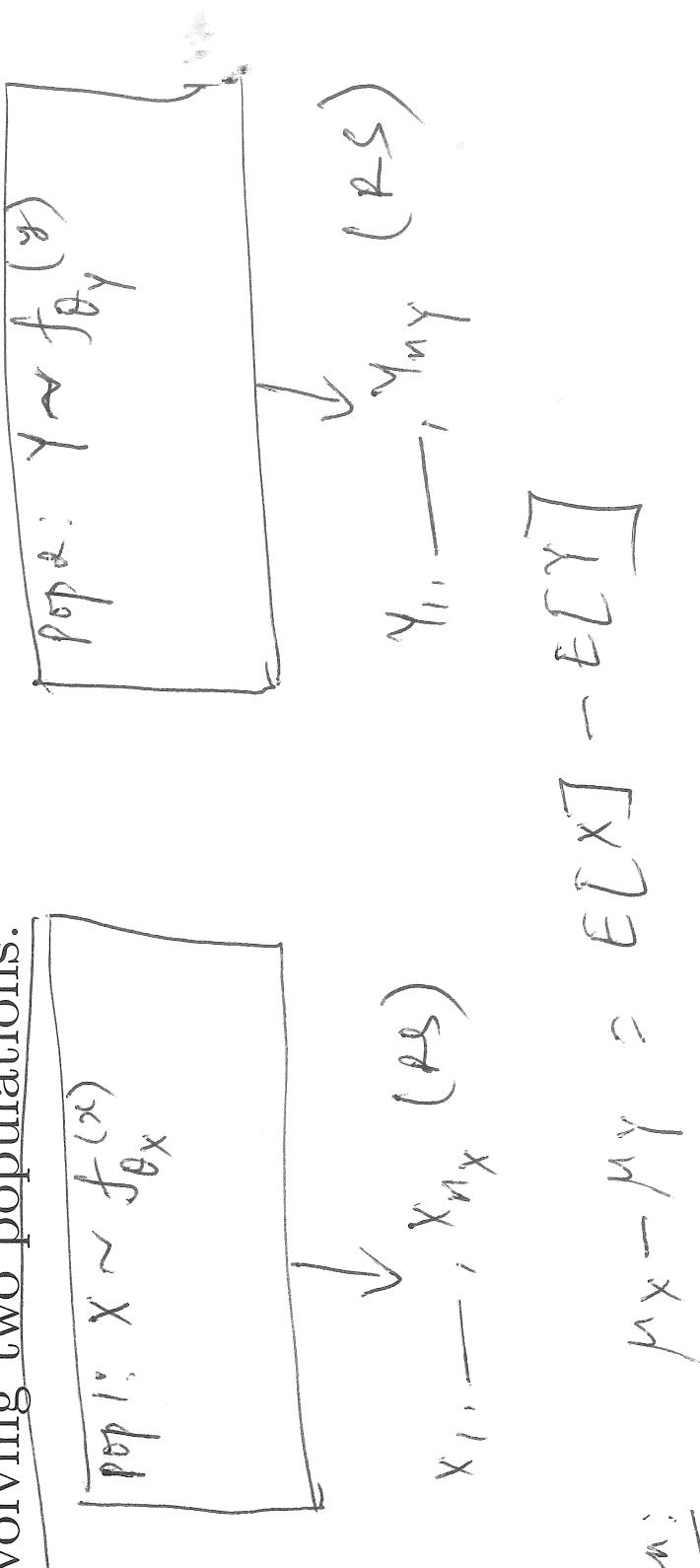
so far:

one-sample problems

Two-sample problems

So far: One-sample problems — inference on parameter(s) of a single population.

Now: Two-sample problems — inference on parameters involving two populations.

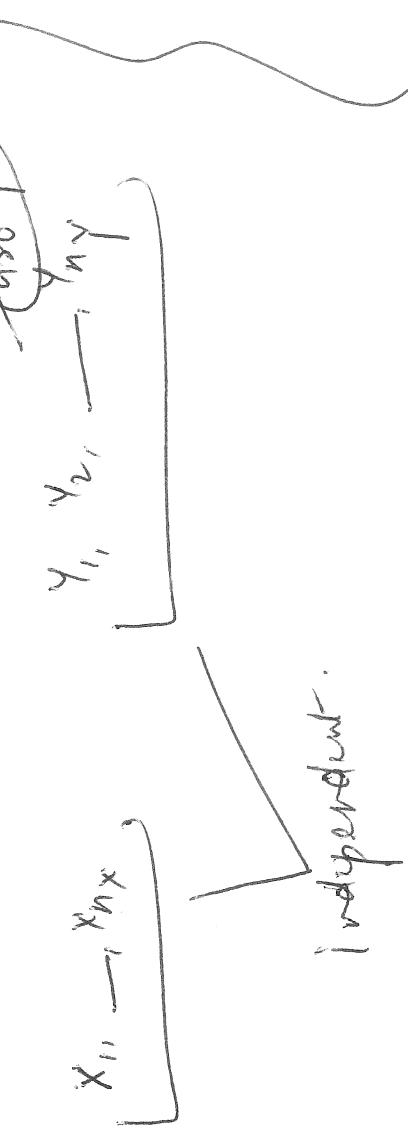


Often:

$$\mu_X - \mu_Y = E[X] - E[Y]$$

Independent or paired samples?

Design 1: (Two independent samples)



Design 2: (Paired samples)

ID	(x_i, y_i)
1	(x_1, y_1)
2	(x_2, y_2)
⋮	⋮
n	(x_n, y_n)

Generally:

Assume different subjects are found.

Given a choice, if feasible, a paired design is better than the two-sample (indep) design. This is because in case of the former, any diff. in X and Y is just due to the effect of the intervention in study. In case of the latter, the effect of the intervention gets confounded with the difference in the groups.