

Recap.

• $X \sim N[\mu, \sigma^2]$, σ^2 known.

$$100(1-\alpha)\% \text{ CI for } \mu: \bar{X} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

- pivot: $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0,1)$.

• - sample size formula:

$$n \approx \left\lceil \frac{2 z_{\alpha/2}^2 \sigma^2}{W} \right\rceil \quad (\text{Round up}).$$

$X \sim N(\mu, \sigma^2)$, σ^2 unknown.

Result: CI for μ : $\bar{X} \pm t_{\alpha/2, n-1} S/\sqrt{n}$

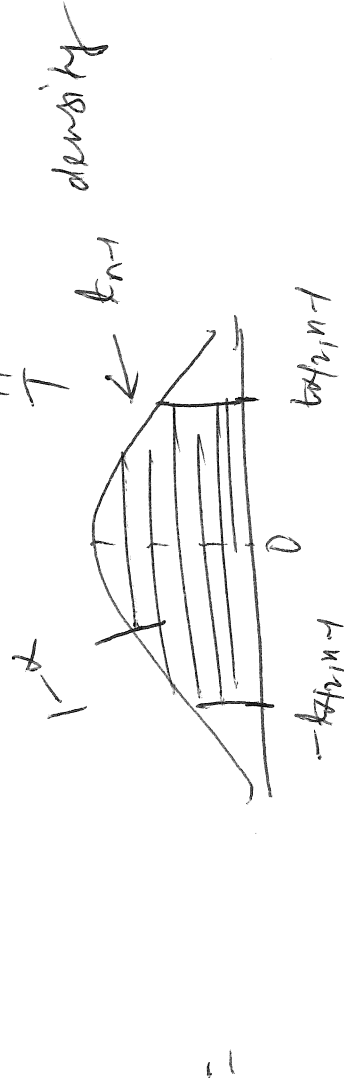
Proof:

Proof: $T = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{n-1}$

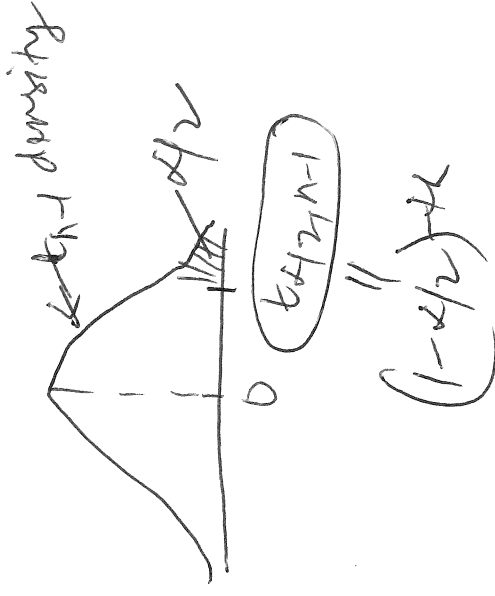
Need to verify that: coverage prob. of this CI = $1 - \alpha$.

$$P\left[\bar{X} - t_{\alpha/2, n-1} \frac{S}{\sqrt{n}} \leq \mu \leq \bar{X} + t_{\alpha/2, n-1} \frac{S}{\sqrt{n}} \right]$$

$$= P\left[-t_{\alpha/2, n-1} \leq \frac{\bar{X} - \mu}{S/\sqrt{n}} \leq t_{\alpha/2, n-1} \right]$$



= $1 - \alpha$ for all (μ, σ^2) .



quantile of t_{n-1} dist.

$$P[t_{n-1} \geq t_{\alpha/2, n-1}] = \alpha/2$$

Note: If n is large: $t_{\alpha/2, n-1} \approx z_{\alpha/2}$

- The t critical points are tabulated in the t -table. Alternatively, we can use `qt` function in R.
- Sample size calculation now becomes complicated than before because S needs to be known before data are collected.
- One option is to make an intelligent guess about S and be conservative (guess a larger value of S so that n larger than necessary is chosen).

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:

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.

$$\begin{array}{l}
 X \sim N[\mu, \sigma^2] \xrightarrow{\text{unknown}} t\text{-interval.} \\
 \downarrow \\
 \text{time b/w keystroke for user} \quad \bar{X} \pm t_{\alpha/2, n-1} \cdot \frac{s}{\sqrt{n}} \\
 \text{trying to log in}
 \end{array}$$

```
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  $\bar{x}$ 
```

```
#> sd(x)
```

```
# [1] 0.1183216  $s$ 
```

```
#> qt(0.975, 11)  $\rightarrow d.f.$ 
```

```
# [1] 2.200985
```

1-42 for 95% confidence

$$95\% \text{ CI: } 0.3 \pm 2.2 \frac{0.118}{\sqrt{12}}$$

$$= 0.3 \pm 0.07$$

$$= [0.23, 0.37] s.$$

$$\text{suppose } \mu_{\text{account owner}} = 0.4 \text{ sec}$$

plausible values of μ .

(ppt) mean time b/w
key strokes for
user trying to log in.

$$\begin{array}{c} 1 \\ 0 \end{array} [0.23 \quad 0.37] \leftarrow$$

II. I think of this as
a constraint to the

α , the right-tail probability	(d.f.)
.0001	1
.0005	2
.001	3
.0025	4
.005	5
.01	6
	7
	8
	9
	10
	11
	12
	13
	14
	15
	16
	17
	18
	19
	20
	21
	22
	23
	24
	25
	26
	27
	28
	29
	30
	32
	34
	36
	38
	40
	45
	50
	55
	60
	70
	80
	90
	100
	200
	∞

Hint: $X \sim \text{some distribution with mean } \mu$.

Large sample CI for mean μ

Recall: When n is large, an approximate $100(1 - \alpha)\%$ CI for mean μ of any population is : $\bar{X} \pm Z_{\alpha/2} \frac{s}{\sqrt{n}}$

Pivot: $Z = \frac{\bar{X} - \mu}{s/\sqrt{n}} \sim N(0,1)$ if n is large.

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 - don't know the distn. but $n=35 \Rightarrow$ large-sample CI.

"
Execution time of the program

$$\bar{X} \pm Z_{\alpha/2} \frac{s}{\sqrt{n}} = 230 \pm 1.96 \frac{14}{\sqrt{35}}$$

$$= \left[\underline{\hspace{1cm}}, \overline{\hspace{1cm}} \right] \text{ ms.}$$

Large sample CI for success proportion p

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

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

Recall: Estimator for $p = \hat{p} = \text{proportion of successes in the sample}$.

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

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

~~Proof:~~

Recall:

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

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

\Rightarrow

large-sample interval:

$$\bar{X} \pm z_{\alpha/2} \hat{SE}(\bar{X})$$

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.

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

↑
Indicator of chip performance ('1' or '0').

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

95% CI for p :

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

$$= 0.4 \pm 1.96 \sqrt{\frac{(0.4)(0.6)}{50}}$$

$$= [0.28, 0.52]$$

Choosing the sample size n

- Width of CI = $2 z_{\alpha/2} \sqrt{\hat{p}(1-\hat{p})} = w$
- 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

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

$\nearrow \hat{p}$ is unknown. ← same thing "worst case"

- 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.)

$$\Rightarrow n \geq \left[\frac{2 z_{\alpha/2} \sqrt{\frac{1}{4}}}{w} \right]^2 = \left[\frac{z_{\alpha/2}}{w} \right]^2$$

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

$$n \geq \frac{1-\alpha}{\left[\frac{z_{\alpha/2}}{w} \right]^2} = \left[\frac{1.96}{0.06} \right]^2 \approx 1068$$

width

So far:

One-sample problem.

Now:

Two-sample problem.

$$X \sim f_{\theta_1}(x)$$

~~x_1, x_2, \dots, x_n~~
 x_1, x_2, \dots, x_{n_1}

$$Y \sim f_{\theta_2}(x)$$

x_1, x_2, \dots, x_{n_2}

Are these samples "independent" or "paired"?

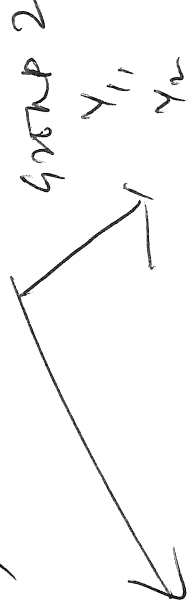
First thing:

Two designs:

Design 1: (Two indep. samples)

The two groups have different subjects.

indep. subjects



Group 1

x_1, x_2, \dots, x_{n_1}

Group 2

y_1, y_2, \dots, y_{n_2}

Design 2:

paired design

Subject	X	Y
1	x_1	y_1
2	x_2	y_2
...		
n	x_n	y_n

"paired data"
~~the~~ X and
the Y are
samples
not
indep.