

# STAT 515 fa 2023 Lec 10 slides

## Confidence intervals for the mean and proportion

Karl B. Gregory

University of South Carolina

Wed: Up to 2  
bonus points  
on exam.

Quiz

These slides are an instructional aid; their sole purpose is to display, during the lecture, definitions, plots, results, etc. which take too much time to write by hand on the blackboard. They are not intended to explain or expound on any material.

SC poll from Sep, 2020: From a sample of 824 SC voters, 47% of them said they would vote for Trump, 43% for Biden, 1% for Jo Jorgensen, 1% for Howie Hawkins, and 8% are not sure.

$$95\%: \hat{p}_n \pm 1.96 \sqrt{\frac{\hat{p}_n(1-\hat{p}_n)}{n}} = .47 \pm 1.96 \sqrt{\frac{.47(1-.47)}{824}} = 0.034$$

Margin of error reported as 3.4%, so if  $p$  is the proportion for Trump, poll says

$$p \in \underbrace{(0.47 - 0.034, 0.47 + 0.034)}_{\substack{L \quad U}} = (0.436, 0.504).$$

- This type of interval is called a *confidence interval (CI)*.
- We like to calibrate CIs so they capture their target with probability  $1 - \alpha$ .
- The value  $\alpha \in (0, 1)$  is called the *confidence level*.

↑ alpha

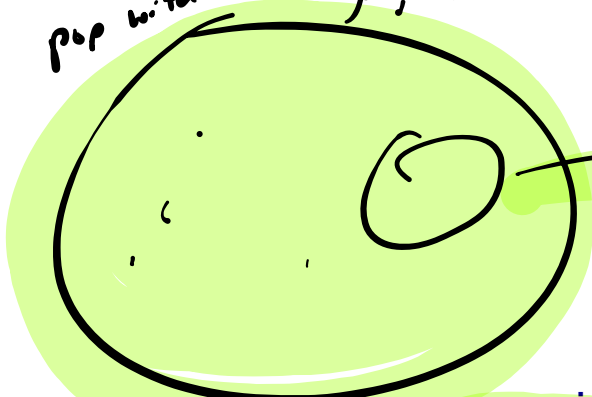
↑  
 $\alpha$  is how often we permit ourselves to miss the target.

The idea of CIs is to find lower and upper bounds  $L$  and  $U$  such that

$$p \in (L, U) \quad \text{or} \quad \boxed{\mu \in (L, U)} \quad \text{or} \quad \sigma^2 \in (L, U),$$

for example, with probability  $1 - \alpha$ .

pop with mean  $\mu$ , variance  $\sigma^2$



$X_1, \dots, X_n, \bar{X}_n$  sample.

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

Want to make an interval which "captures"  $\mu$  w/ prob. 0.95.

$$\bar{X}_n \pm ME$$

Exercise: Let  $X_1, \dots, X_n \stackrel{\text{ind}}{\sim} \text{Normal}(\mu, \sigma^2)$  and use the fact that

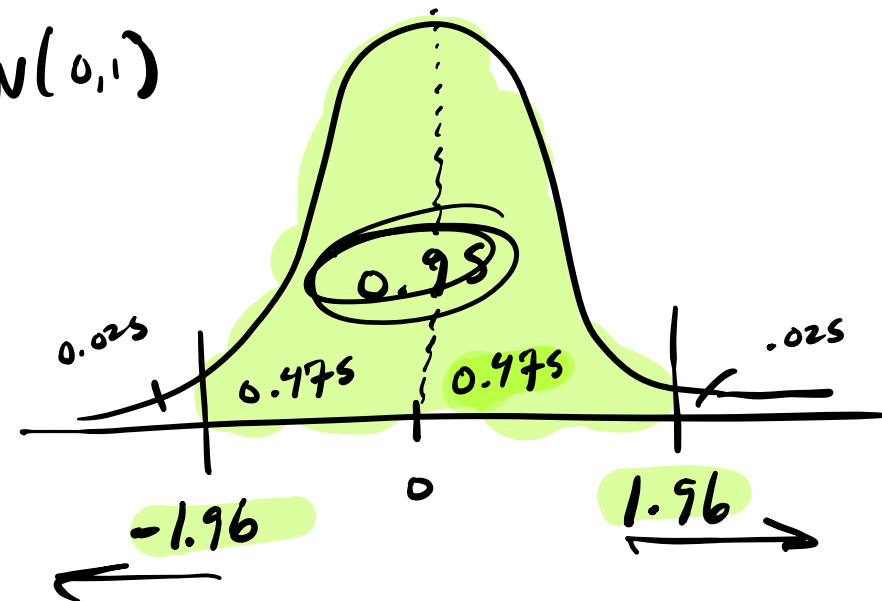
$$P\left(-1.96 \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq 1.96\right) = 0.95$$

to find a 95% confidence interval for  $\mu$ .

$$P\left(-\frac{\sigma}{\sqrt{n}} 1.96 < \bar{X} - \mu < \frac{\sigma}{\sqrt{n}} 1.96\right) = 0.95 \quad Z \sim N(0,1)$$

$$\Leftrightarrow P\left(-\bar{X}_n - \frac{\sigma}{\sqrt{n}} 1.96 < -\mu < -\bar{X}_n + \frac{\sigma}{\sqrt{n}} 1.96\right) = 0.95$$

$$\Leftrightarrow P\left(\bar{X}_n + \frac{\sigma}{\sqrt{n}} 1.96 > \mu > \bar{X}_n - \frac{\sigma}{\sqrt{n}} 1.96\right) = 0.95$$



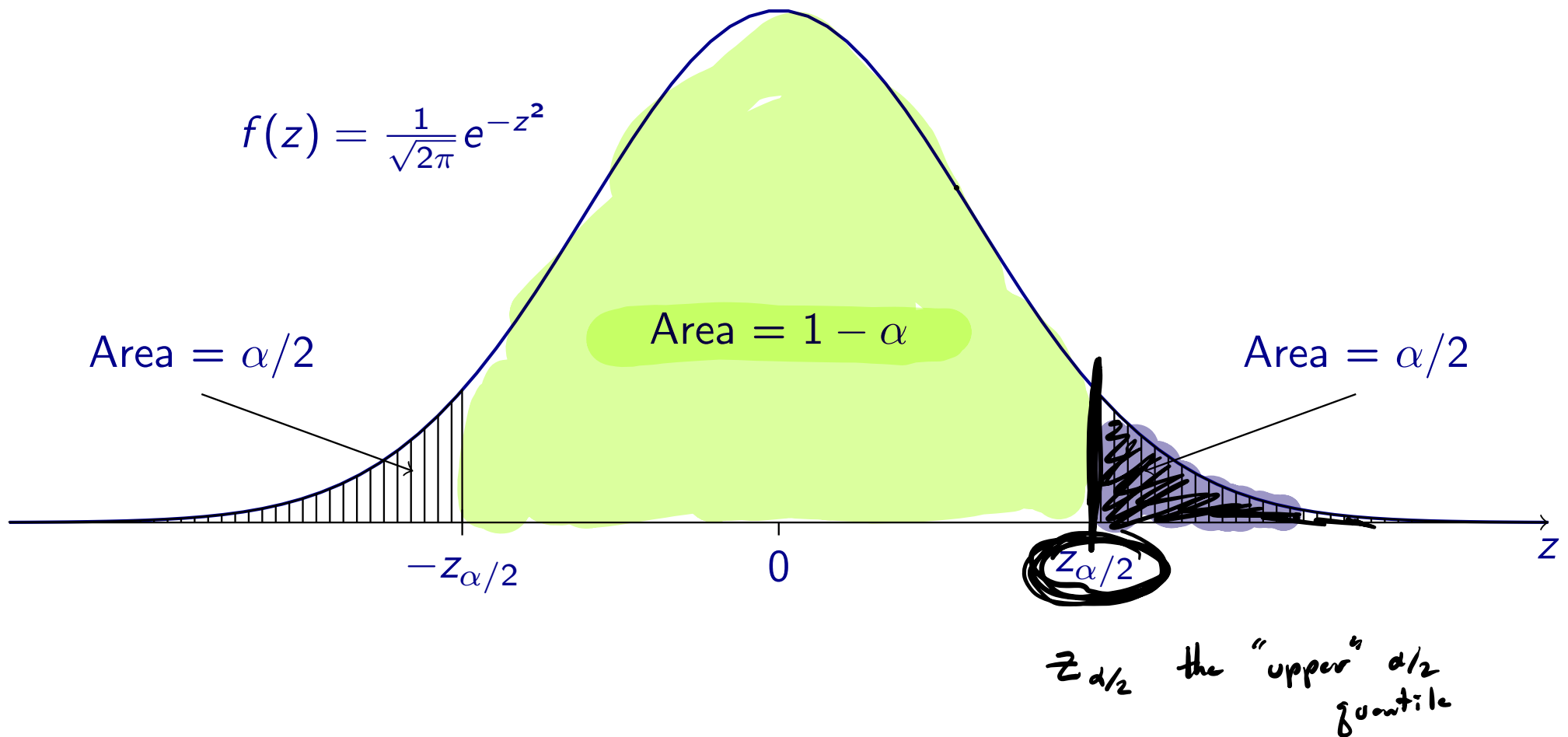
So the interval  $\left( \bar{x}_n - \frac{\sigma}{\sqrt{n}} 1.96, \bar{x}_n + \frac{\sigma}{\sqrt{n}} 1.96 \right)$   
contains  $\mu$  with prob  $0.95$ .

A 95% confidence Interval for  $\mu$  is  
$$\bar{x}_n \pm \frac{\sigma}{\sqrt{n}} \underline{1.96}.$$

95% corresponds to  $1 - \alpha = .95$ ,  $\alpha = \underline{0.05}$ .

$$95\% = (1 - \underbrace{0.05}_{\alpha}) 100\%$$

What about a general  $(1 - \alpha) \times 100\%$  CI for any  $\alpha \in (0, 1)$ ?



## Confidence interval for the mean of a Normal population with $\sigma$ known

Let  $X_1, \dots, X_n \stackrel{\text{ind}}{\sim} \text{Normal}(\mu, \sigma^2)$ . Then a  $(1 - \alpha) \times 100\%$  CI for  $\mu$  is

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

comes from z-table

**Example:** These are the commute times (sec) to class of a sample of students.

1832	1440	1620	1362	577	934	928	998	1062	900
1380	913	654	878	172	773	1171	1574	900	900

Assume the population is Normal with  $\sigma = 400$ .

$X \leftarrow c(1832, 1440, \dots, 900)$   
mean(X)

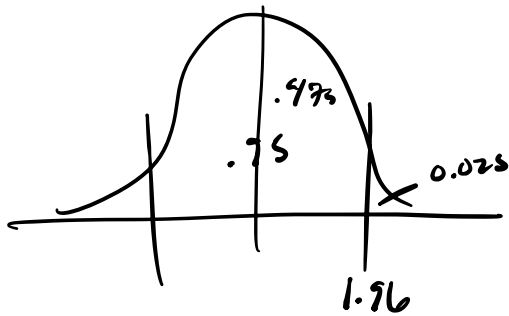
- 1 Construct a 95% confidence interval for the mean commute time of all students.
- 2 Construct a 99% confidence interval for the mean commute time of all students.
- 3 Give an interpretation of the two confidence intervals.
- 4 Which confidence interval is wider? Does it make sense why??

①  $\bar{X}_n = 1048.4$  ,  $\sigma = 400$   $n = 20$

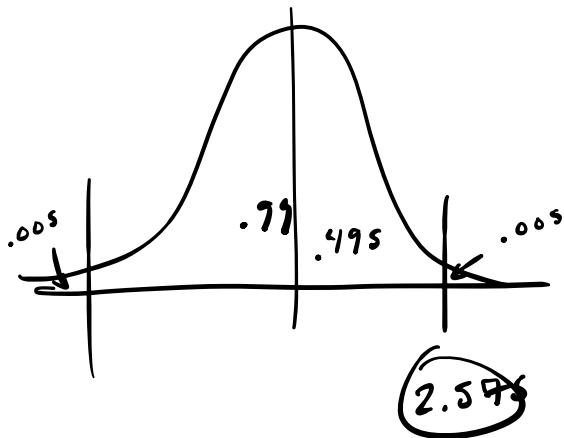
$(1-\alpha) 100\%$  C.I. is  $\bar{X}_n \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$

So a 95% C.I.

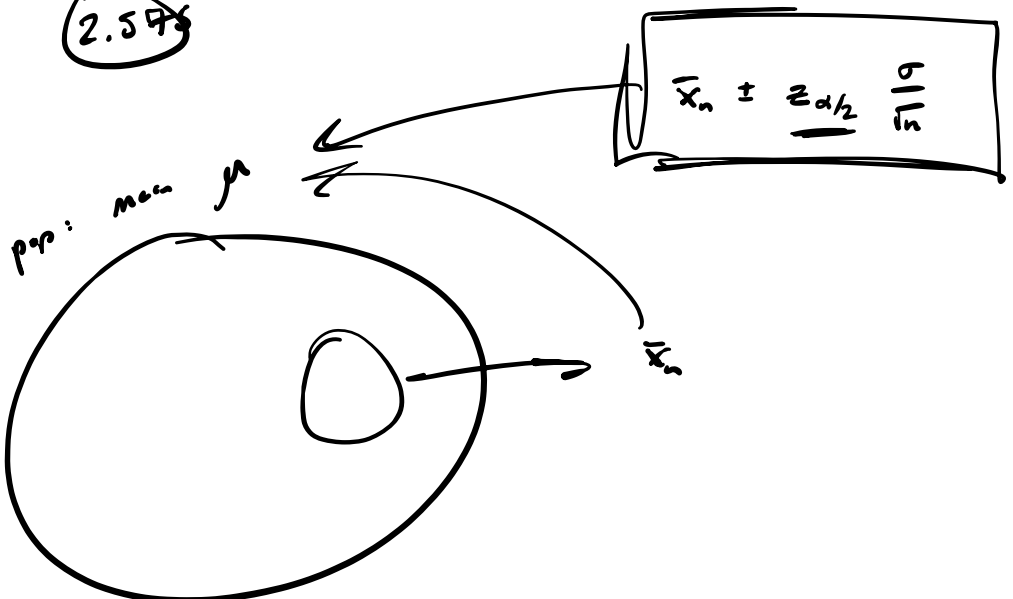
$1048.4 \pm 1.96 \frac{400}{\sqrt{20}} = (873.09, 1223.7)$



② 99% C.I. now:  $\alpha = 0.01$



$1048.4 \pm 2.575 \frac{400}{\sqrt{20}} = (818.1, 1278.7)$



## Confidence interval for mean of a non-Normal pop. with $\sigma$ known

Let  $X_1, \dots, X_n$  be a rs from a non-Normal dist. with mean  $\mu$  and var.  $\sigma^2 < \infty$ .  
Then

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

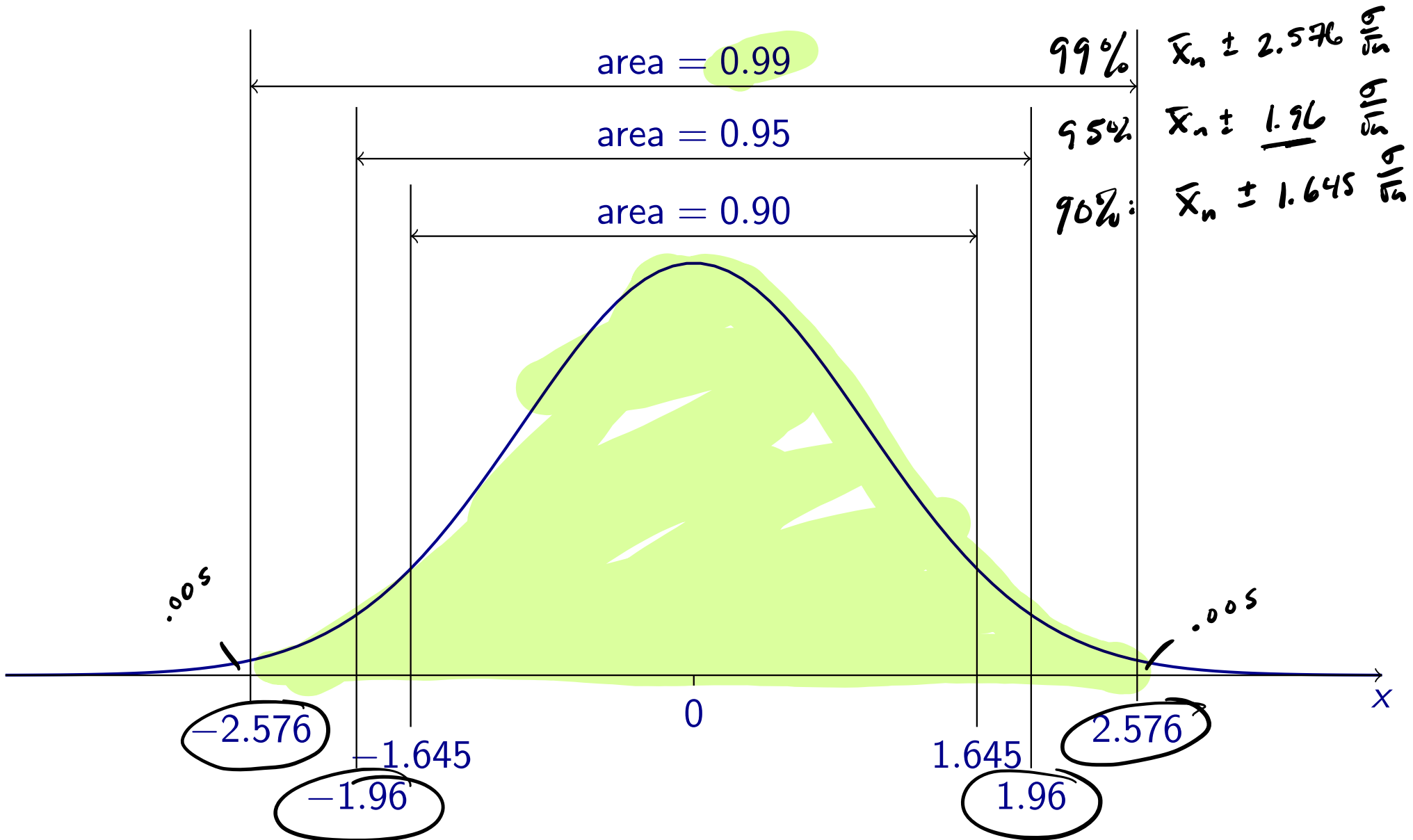
contains  $\mu$  with probability closer and closer to  $1 - \alpha$  for larger and larger  $n$ .

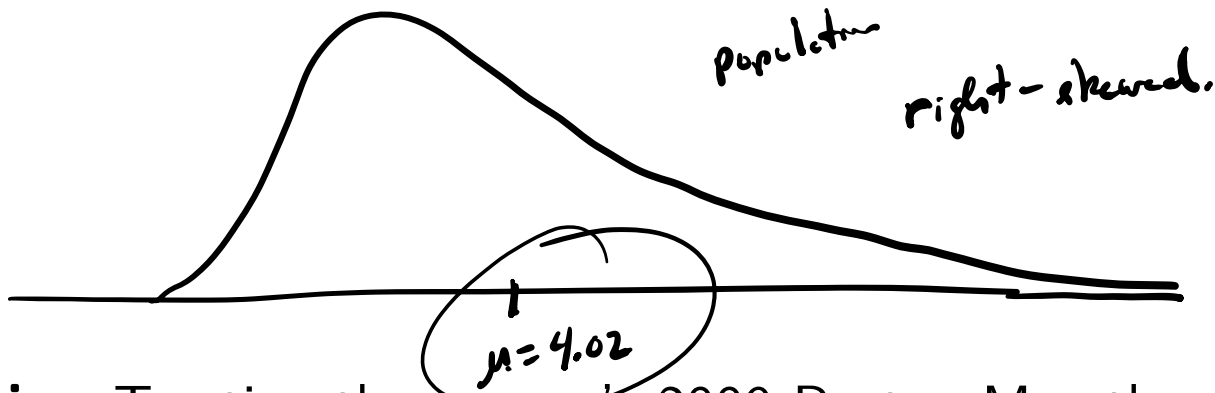
This is a **Central Limit Theorem** result.

$$Z = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \stackrel{\text{approx}}{\sim} \text{Normal}(0, 1) \quad \text{for large } n$$



We often construct 90%, 95%, and 99% confidence intervals:





$$\sigma = 0.555$$

**Exercise:** Treating the women's 2009 Boston Marathon finishing times, which have mean  $\mu = 4.02$  and standard deviation  $\sigma = 0.555$ ,

- 1 Draw 100 samples of size  $n = 30$  and each time build a 90% CI for the mean.
- 2 Check for what proportion of samples the CI captured the true mean  $\mu$ .

[link to women's data.](#)

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

Draw sample of  $n = 30$  finishing times.

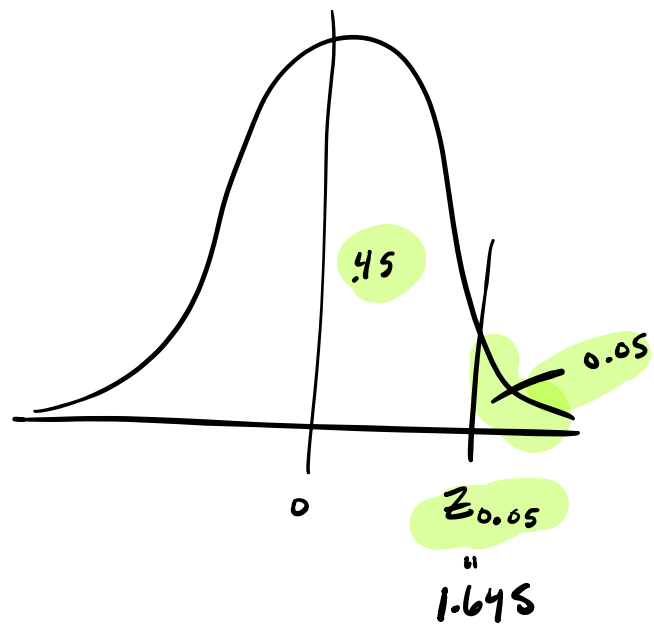
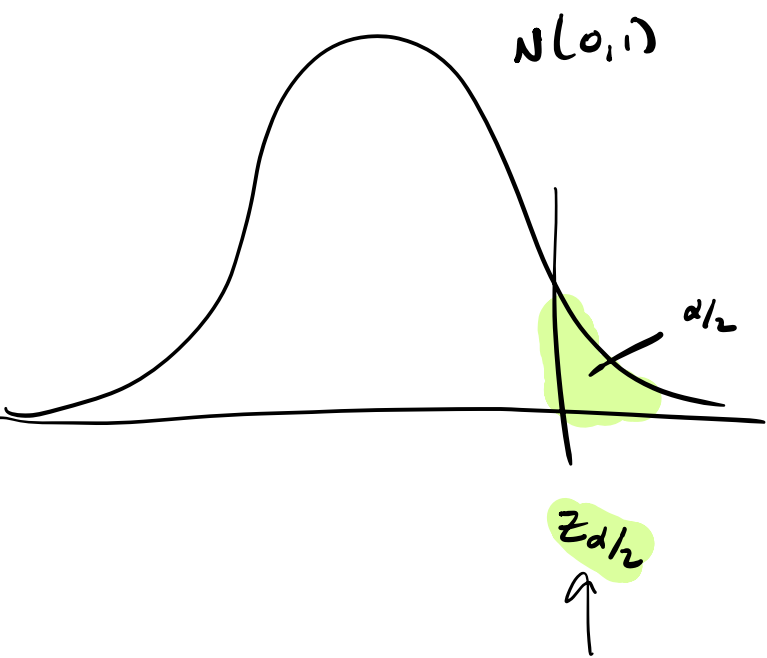
90% C.I.

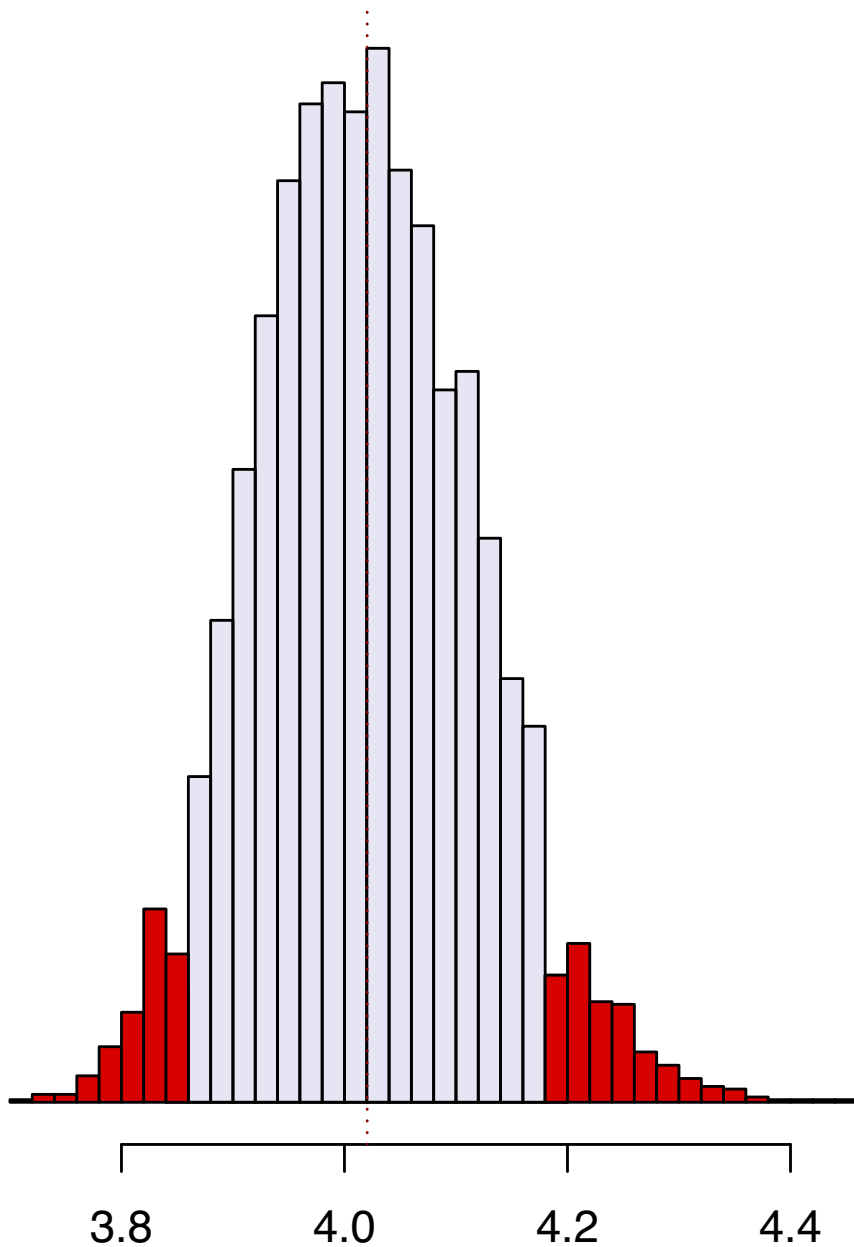
$$\bar{X}_n \pm z_{\frac{0.1}{2}} \frac{0.555}{\sqrt{30}}$$

$$\alpha = 0.10$$

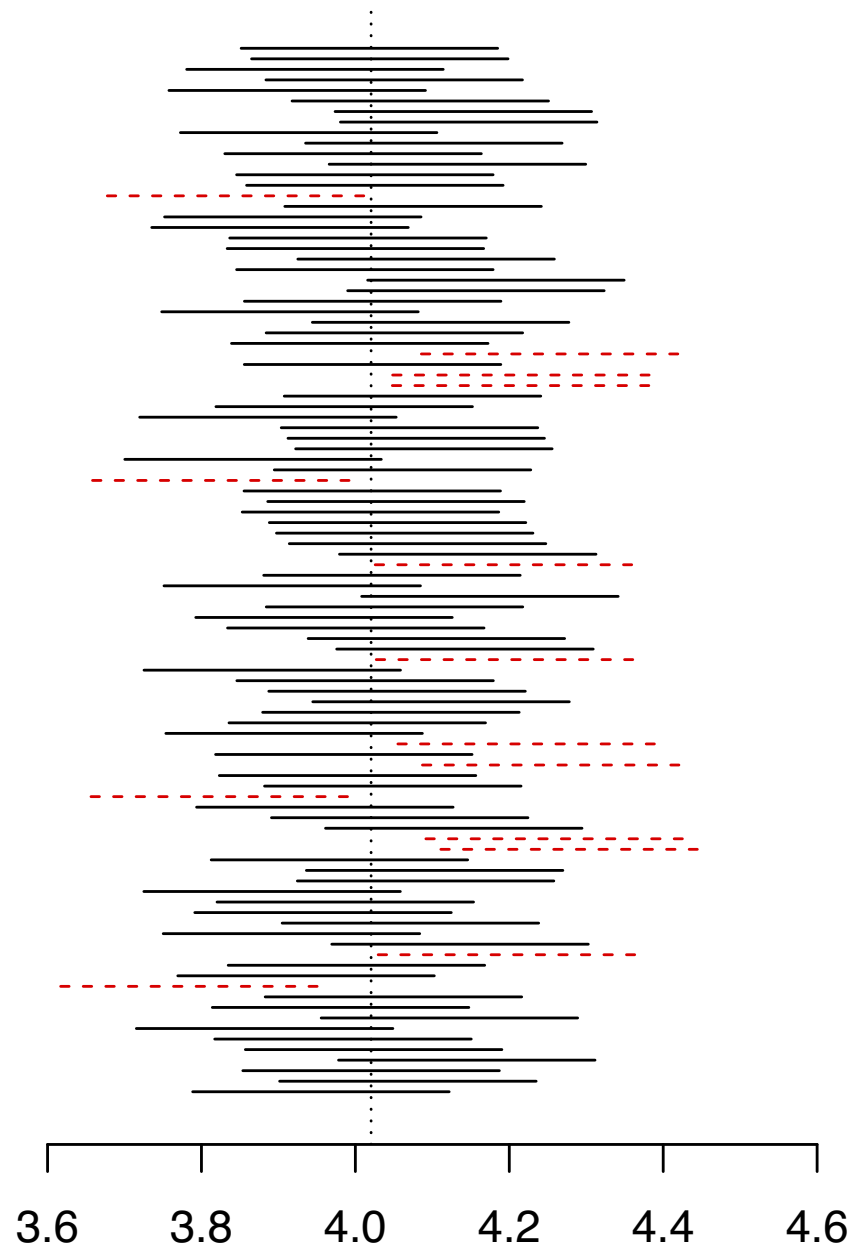
$$z_{.05} = 1.645$$

$$\bar{X}_n \pm 1.645 \frac{0.555}{\sqrt{30}}$$





$\bar{X}_n$

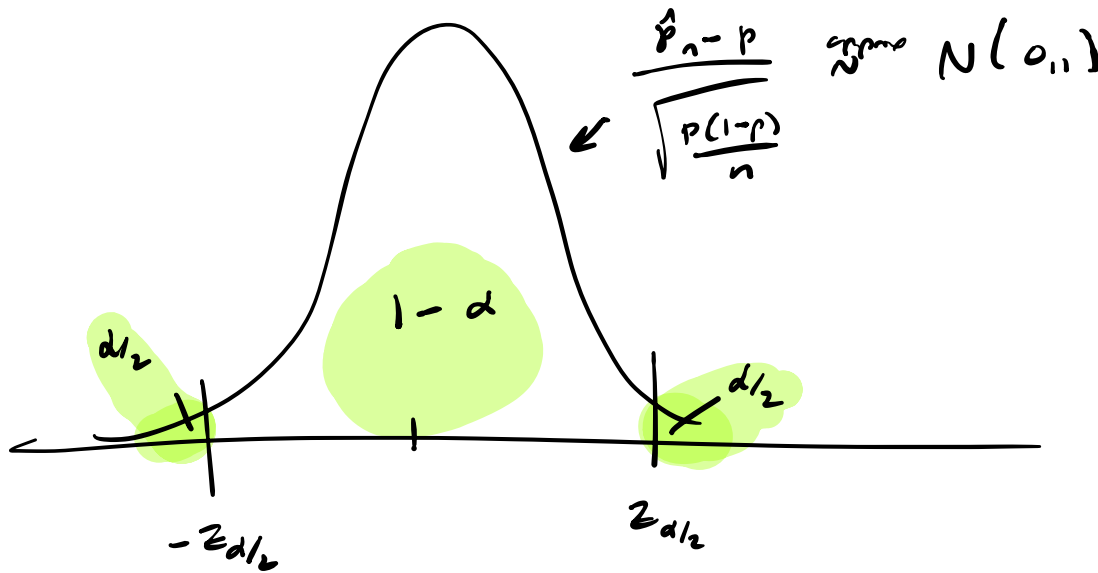


$(\bar{X}_n - 1.645 \cdot \sigma/\sqrt{n}, \bar{X}_n + 1.645 \cdot \sigma/\sqrt{n})$

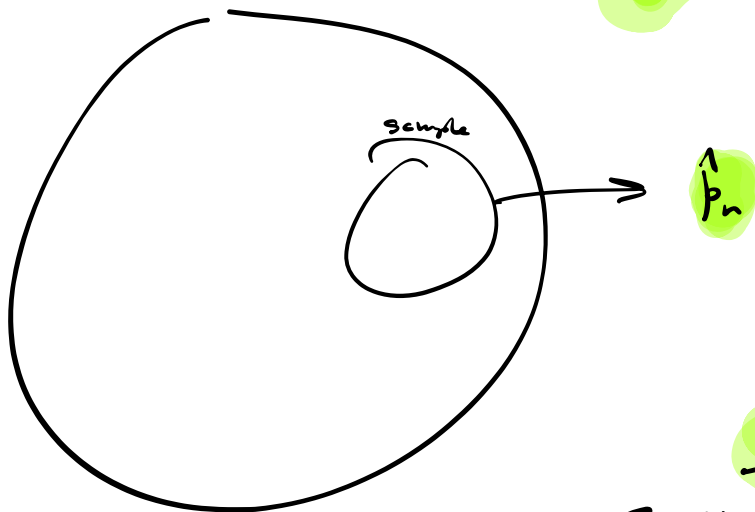
rearrange this

$$P\left(-z_{d/2} < \frac{\hat{p}_n - p}{\sqrt{\frac{p(1-p)}{n}}} < z_{d/2}\right) = 1 - \alpha$$

$$P\left(\hat{p}_n - z_{d/2} \sqrt{\frac{p(1-p)}{n}} < p < \hat{p}_n + z_{d/2} \sqrt{\frac{p(1-p)}{n}}\right) = 1 - \alpha$$



pop: proportion of 1s is  $p$ .



$$Z \sim \frac{\hat{p}_n - p}{\sqrt{\frac{p(1-p)}{n}}} \stackrel{\text{approx}}{\sim} N(0,1)$$

if n large.

$$E \hat{p}_n = p$$

$$\text{Var } \hat{p}_n = \frac{p(1-p)}{n}$$

pl

## Confidence interval (Wald-type) for a proportion

If  $n\hat{p}_n \geq 15$  and  $n(1 - \hat{p}_n) \geq 15$ , then

# successes

# failures

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

is an approximate  $(1 - \alpha) \times 100\%$  CI for  $p$ .

Explain how we arrive at this confidence interval.

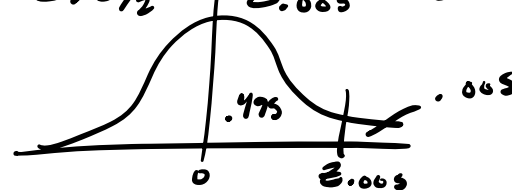
**Exercise:** From a sample of 1,000 voters, 478 say they will vote for candidate A. Let  $p$  be the true proportion of voters who will vote for candidate A.

① Build a 95% CI for  $p$ .  $\alpha = 0.05 \rightarrow z_{\alpha/2} = z_{0.05/2} = z_{0.025} = 1.96$

② Build a 99% CI for  $p$ .  $\alpha = 0.01 \rightarrow z_{\alpha/2} = z_{0.01/2} = z_{0.005} = 2.575$

①  $\hat{p}_n = \frac{478}{1000}$

$n = 1000$



$$\hat{p}_n \pm z_{d/2} \sqrt{\frac{\hat{p}_n(1-\hat{p}_n)}{n}} = \frac{478}{1000} \pm 1.96 \sqrt{\frac{\frac{478}{1000} \left(1 - \frac{478}{1000}\right)}{1000}}$$

$$= .478 \pm 0.031$$

$$= (0.447, 0.509)$$

$$47.8\% \pm 3.1\%$$

$$= \frac{478}{1000} \pm 2.576 \sqrt{\frac{\frac{478}{1000} \left(1 - \frac{478}{1000}\right)}{1000}} = (0.434, 0.519)$$

Sample of size  $n$

$$\hat{p}_n = \frac{\# \text{successes}}{n}$$

## Confidence interval (Agresti-Coull) for a proportion

If  $n\hat{p}_n \geq 5$  and  $n(1 - \hat{p}_n) \geq 5$ , then

$$\hat{p}_n \pm z_{\alpha/2} \sqrt{\frac{\tilde{p}_n(1 - \tilde{p}_n)}{n + 4}}, \quad \text{where } \tilde{p}_n = \frac{\# \{ \text{successes} \} + 2}{n + 4}$$

is an approximate  $(1 - \alpha) \times 100\%$  CI for  $p$ .

$$n = 50 \quad \hat{p}_n = \frac{5}{50} \quad \tilde{p}_n = \frac{7}{54}$$

✓ **Exercise:** From a sample of 50 students, 5 say they hang-dry their laundry. Let  $p$  be the true proportion of students who hang-dry their laundry.

- ① Build a 95% Wald-type CI for  $p$ .  $\hat{p}_n \pm 1.96 \sqrt{\frac{\hat{p}_n(1 - \hat{p}_n)}{n}} = (0.017, 0.183)$
- ② Build a 95% Agresti-Coull CI for  $p$ .  $\tilde{p}_n \pm 1.96 \sqrt{\frac{\tilde{p}_n(1 - \tilde{p}_n)}{n + 4}} = (0.04, 0.219)$
- ③ Do the same if 50 out of 500 students say they hang-dry their laundry.

Return to 2020 poll on first slide and check what  $\alpha$  is...