

MODULE 9

INDEPENDENT-SAMPLES T-TEST

OUTLINE

- 1 Between-group designs
- 2 Quick review
- 3 Significance testing steps
- 4 Statistical assumptions
- 5 Study questions

OUTLINE

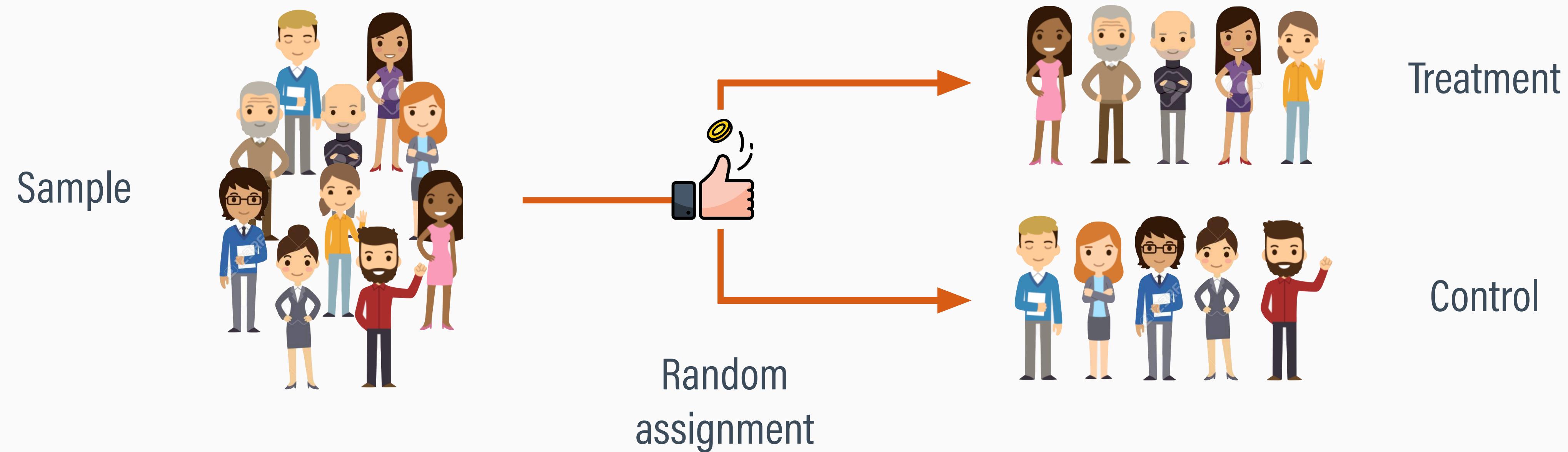
- 1 Between-group designs
- 2 Quick review
- 3 Significance testing steps
- 4 Statistical assumptions
- 5 Study questions

BETWEEN-GROUP RESEARCH DESIGNS

- A **between-group research design** seeks to compare two or more groups of participants
- Unlike the within-group design, each condition is comprised of different participants
- The classic example is a randomized experiment with a treatment and control group, but groups can reflect any qualitative characteristic (e.g., sociodemographic)

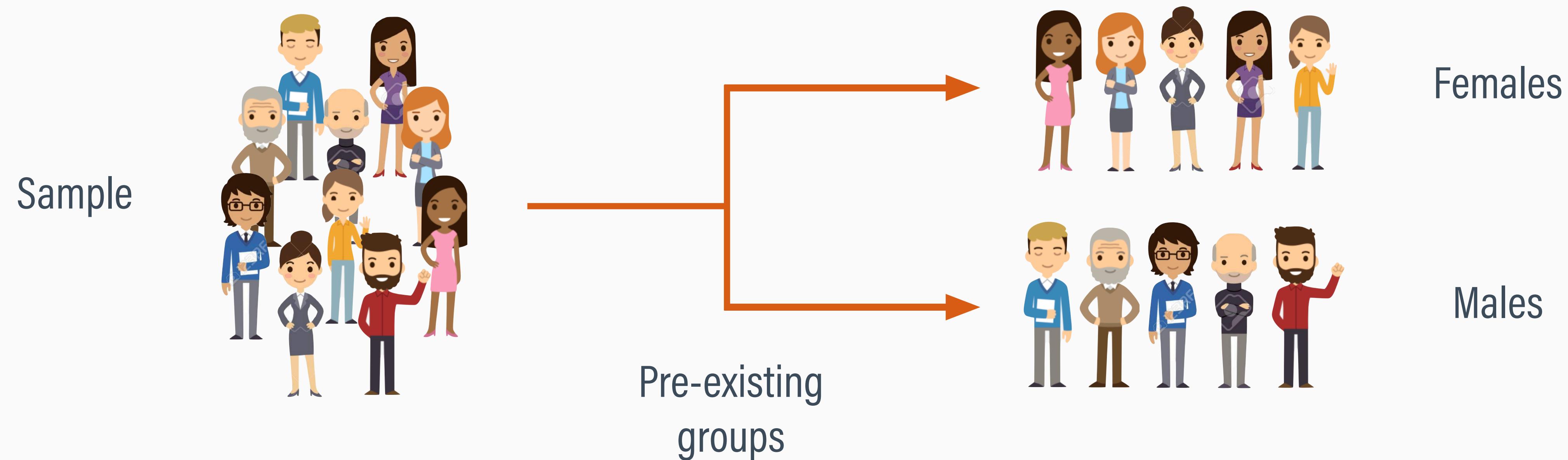
EXPERIMENTAL APPLICATION

- Participants are randomly assigned to either a treatment or a control condition



NON-EXPERIMENTAL APPLICATION

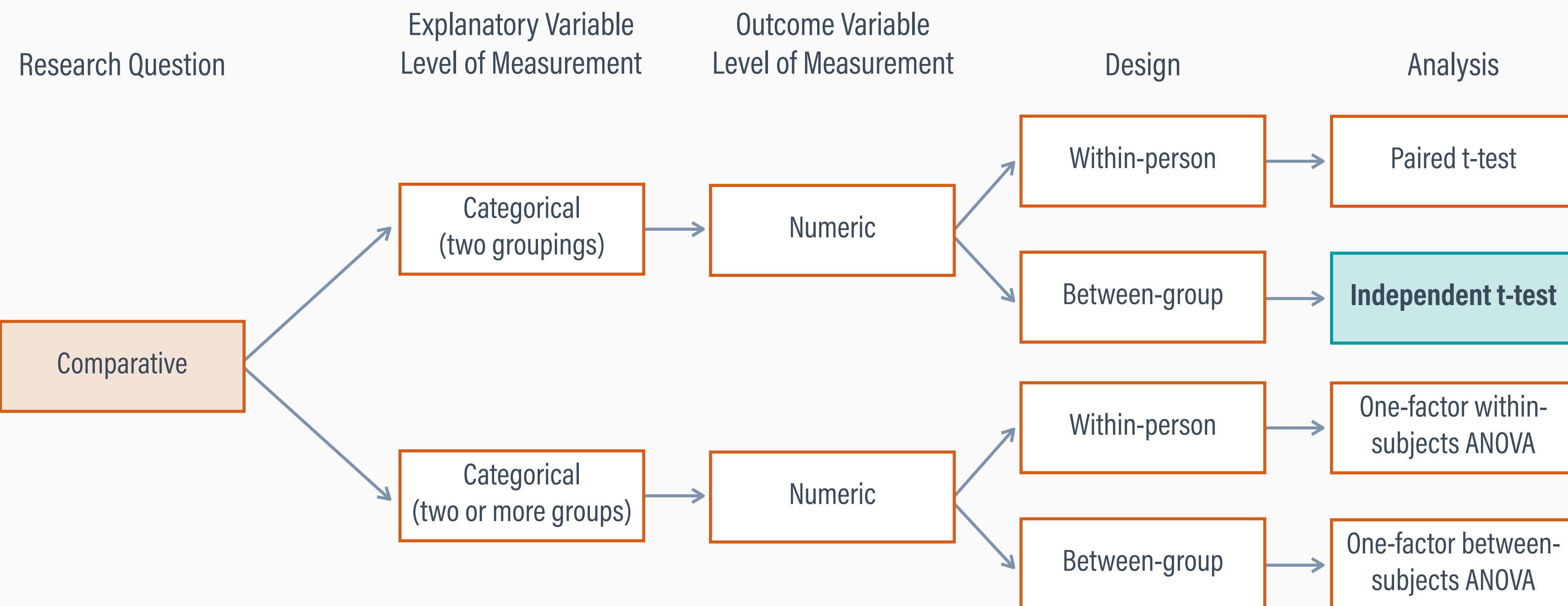
- Participants divide into subgroups based on a shared qualitative characteristic



INDEPENDENT-SAMPLES t TEST

- The independent-samples t-test is appropriate for between-group designs with two groups
- Applicable to comparative research questions and hypotheses involving the difference between two means obtained from the different individuals

STATISTICAL ORG CHART



OUTLINE

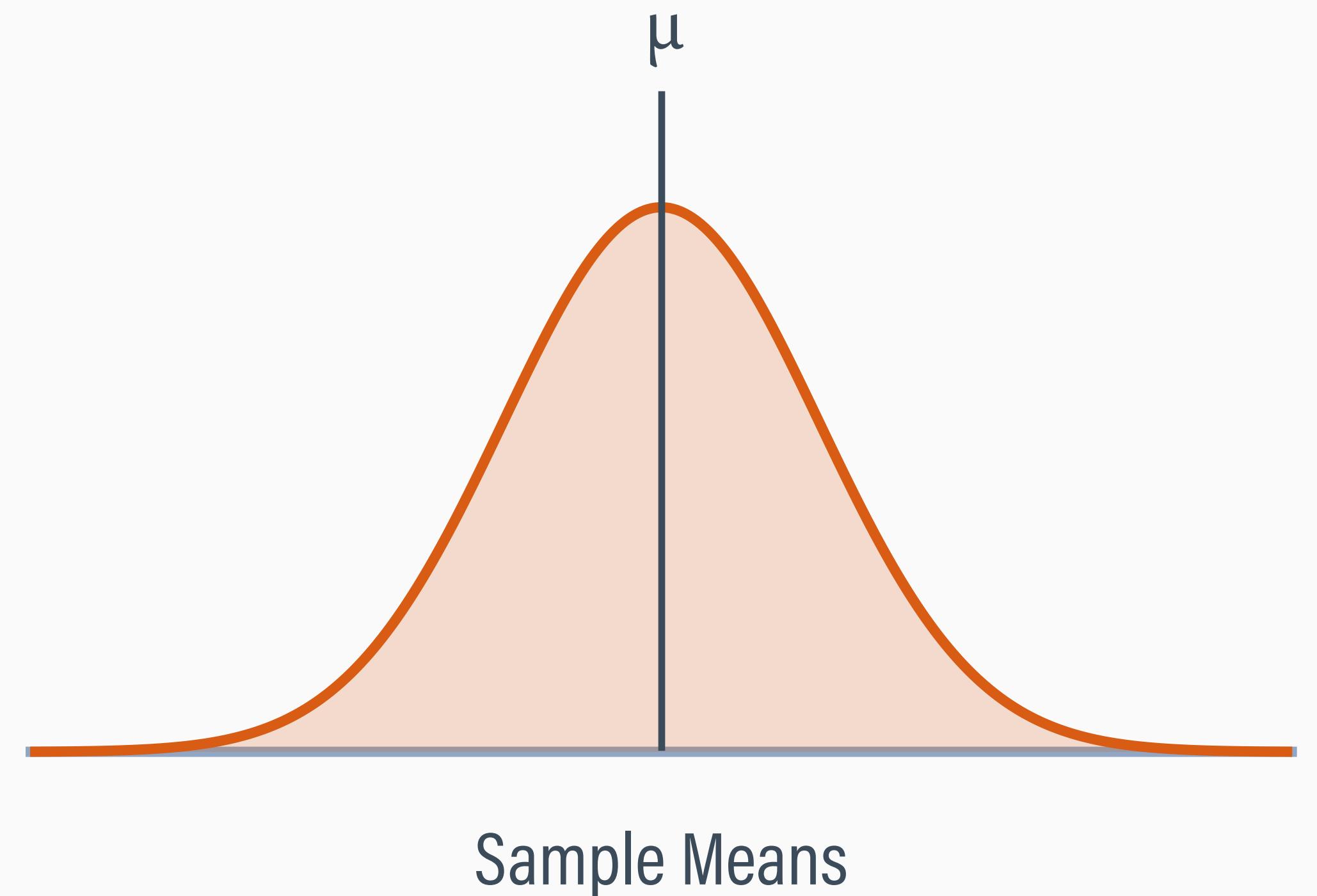
- 1 Between-group designs
- 2 Quick review
- 3 Significance testing steps
- 4 Statistical assumptions
- 5 Study questions

QUICK REVIEW: SAMPLING ERROR

- The frequentist paradigm imagines a single population that spawns many hypothetical random samples of data (one parameter, many hypothetical estimates)
- The amount by which an estimate differs from the true population statistic is called sampling error
- Every hypothetical sample has a different amount of error

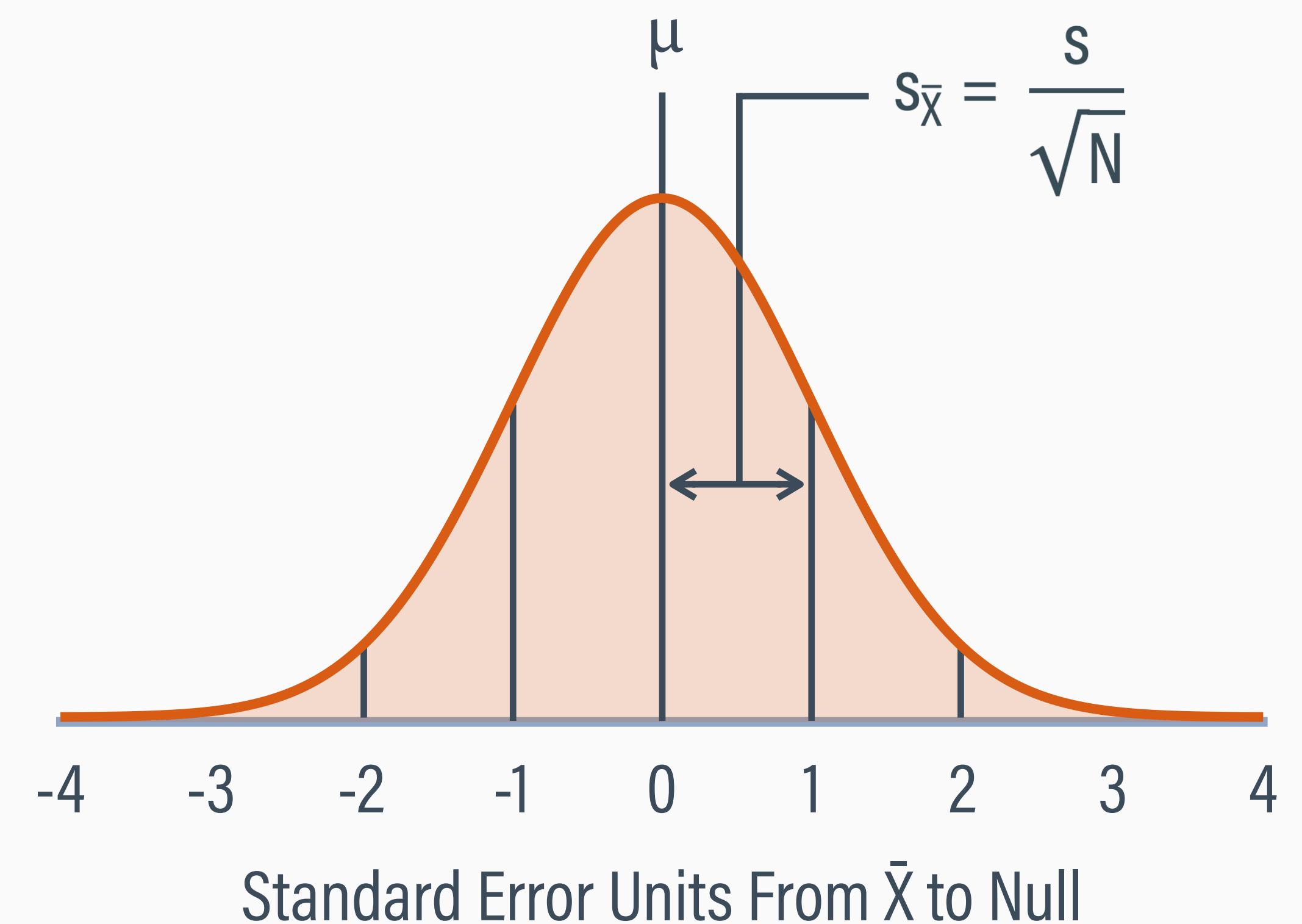
QUICK REVIEW: SAMPLING DISTRIBUTION

- The distribution of the estimates from many hypothetical samples is a sampling distribution
- With a large enough N, sample means follow a normal curve centered at the true mean
- Most estimates have small sampling errors, but a few have larger errors



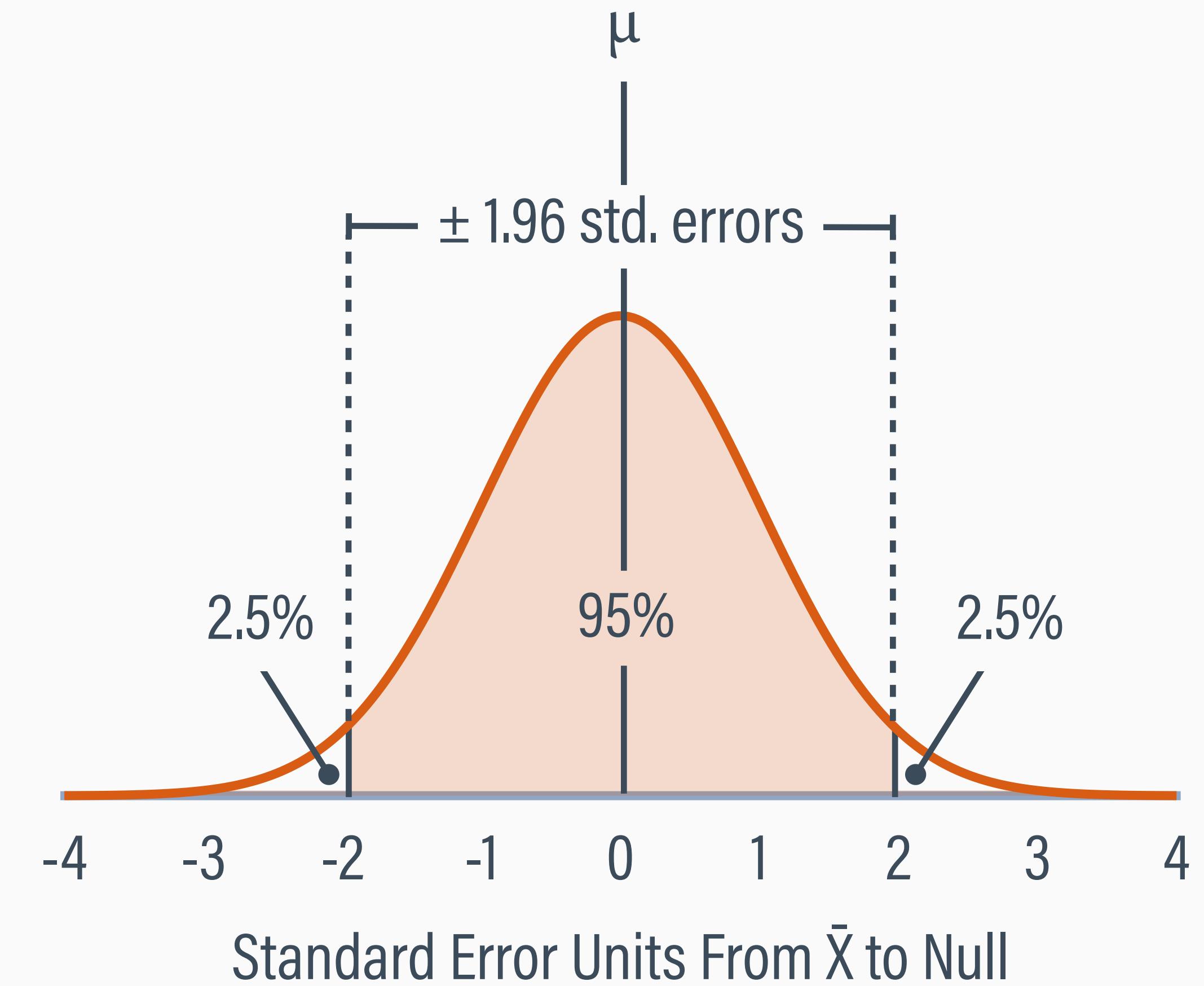
QUICK REVIEW: STANDARD ERROR

- The standard error is the average distance from a sample mean and the true mean
- $s_{\bar{x}} = \text{standard deviation of the sample means}$
- The standard error is the average or expected amount of sampling error across many hypothetical samples



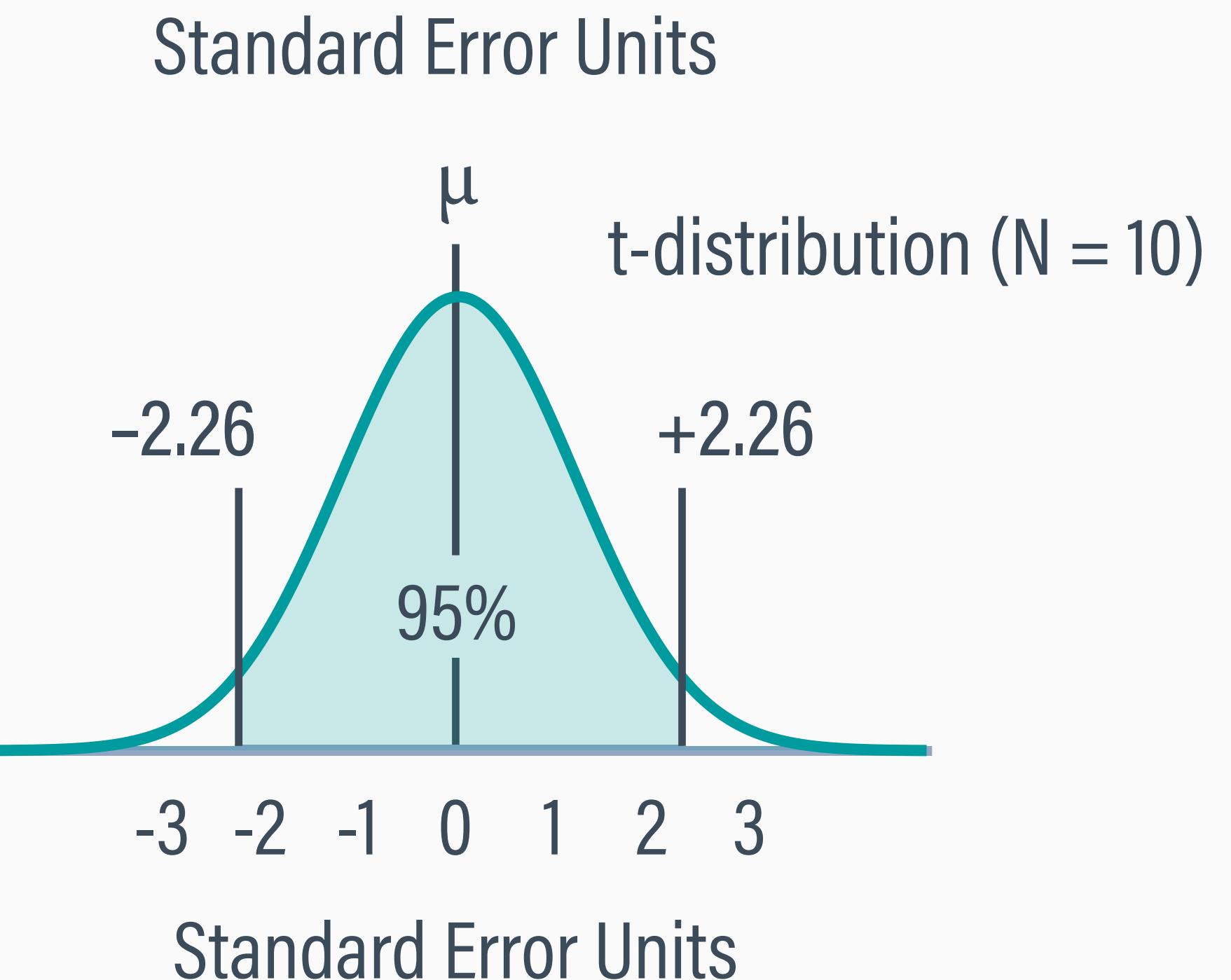
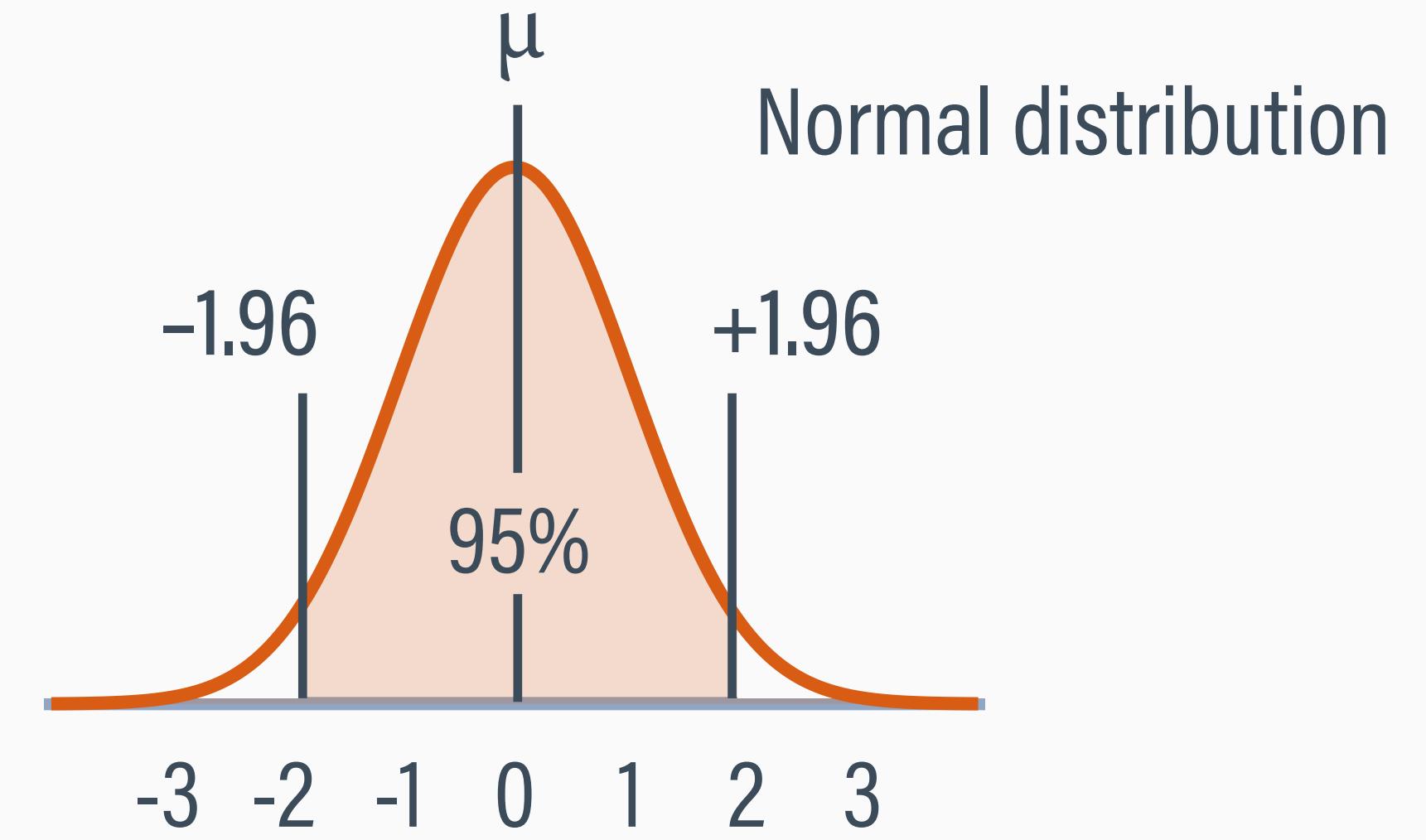
QUICK REVIEW: NORMAL CURVE RULE

- The standard error is the standard deviation of many hypothetical sample means
- We can apply normal curve rules of thumb
- 95% of the means from large samples are within ± 1.96 standard errors of the true mean



QUICK REVIEW: T-DISTRIBUTION

- When using small samples, the normal curve is an inaccurate description of sampling error
- The t-distribution is a series of bell-shaped curves that stretch out (become more variable) as the N decreases
- Small samples are more likely to produce outlier estimates, and “stretching” the curve honors that

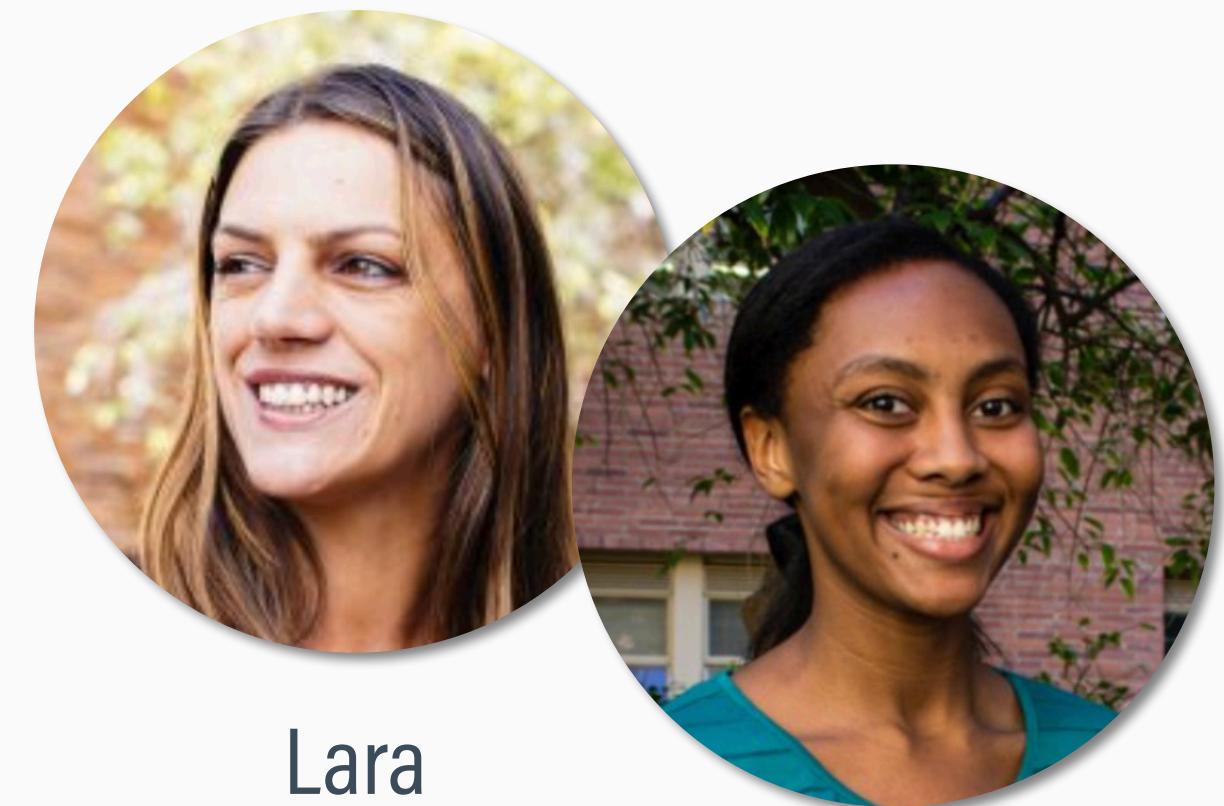


OUTLINE

- 1 Between-group designs
- 2 Quick review
- 3 Significance testing steps
- 4 Statistical assumptions
- 5 Study questions

SMOKING AND DRINKING CESSATION TRIAL

Pharmacological treatments that can concomitantly address cigarette smoking and heavy drinking stand to improve health care delivery for these highly prevalent co-occurring conditions. This superiority trial compared the combination of varenicline and naltrexone against varenicline alone for smoking cessation and drinking reduction among heavy-drinking smokers.

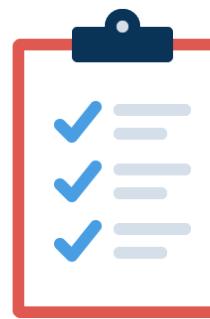


Lara
Ray

ReJoyce
Green

Ray, L.A., Green, R., Enders, C., et al. (2021). Efficacy of combining varenicline and naltrexone for smoking cessation and drinking reduction: A randomized clinical trial. *American Journal of Psychiatry*, 178, 818–828.

KEY VARIABLES



Breath (alveolar) carbon monoxide

A measure of carbon monoxide in the lungs.
Breath carbon monoxide is a biomarker of smoking behavior common in clinical trials.



Medication arm

Participants were randomly assigned to receive one of two meds: varenicline plus naltrexone or varenicline plus placebo pills

RESEARCH QUESTION

- Question: Does smoking intensity differ between participants receiving a combination of two medications and those receiving a single medication alone?
- The explanatory (independent) variable, medication arm, consists of two groups: varenicline only and varenicline plus naltrexone
- The outcome (dependent) variable, breath carbon monoxide, is a numeric biomarker of smoking intensity

SIGNIFICANCE TESTING STEPS

- 1 Specify hypotheses
- 2 Define standard of evidence
- 3 Design study and collect data
- 4 Compare data to null hypothesis
- 5 Evaluate hypotheses and draw conclusion

MEAN DIFFERENCE STATISTIC

- Both experimental groups have population means, μ_1 and μ_2
- Hypotheses use a **mean difference** statistic that contrasts the two population means

$$\mu_{\text{diff}} = \mu_1 - \mu_2$$

- The mean difference quantifies treatment effect in this case

NULL HYPOTHESIS

- In the population, there is no difference between the two medication arms

$$H_0: \mu_{\text{diff}} = 0$$

- The null that $\mu_{\text{diff}} = 0$ is counter to expectations because researchers anticipate that medication differences could differentially impact smoking behavior

TWO POSSIBLE ALTERNATIVE HYPOTHESES

- One-tailed alternate: Smoking intensity could be lower in only one of the medication arms

$$H_A: \mu_{\text{diff}} < 0 (\mu_{V+N} < \mu_V) \quad \text{or} \quad H_A: \mu_{\text{diff}} > 0 (\mu_{V+N} > \mu_V)$$

- Two-tailed alternate: Smoking intensity could be lower in either medication arm

$$H_A: \mu_{\text{diff}} \neq 0 (\mu_{V+N} < \mu_V \text{ or } \mu_{V+N} > \mu_V)$$

SIGNIFICANCE TESTING STEPS

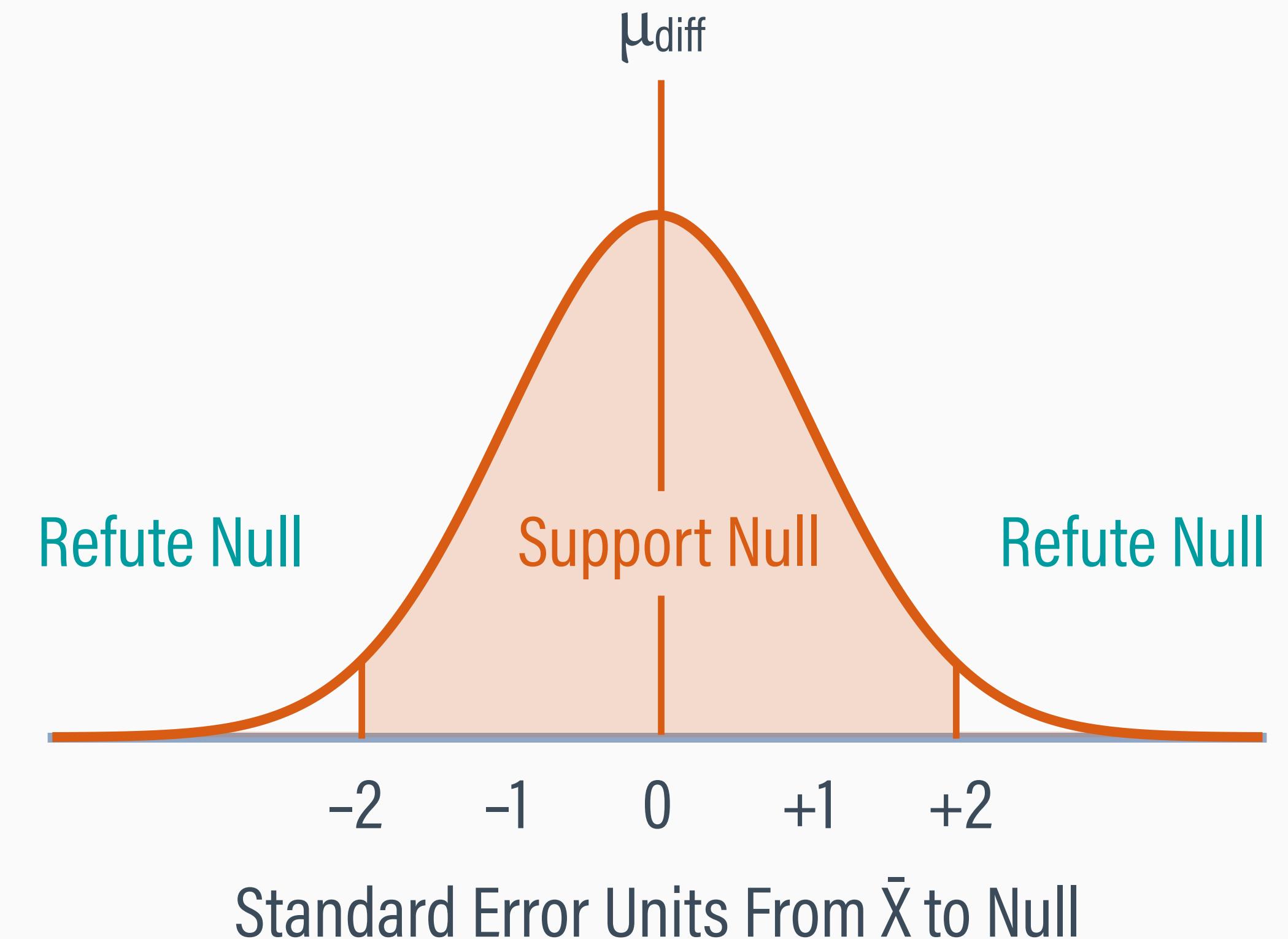
- 1 Specify hypotheses
- 2 Define standard of evidence
- 3 Design study and collect data
- 4 Compare data to null hypothesis
- 5 Evaluate hypotheses and draw conclusion

STANDARD OF EVIDENCE

- The data are the evidence that we use to conclude whether the null is plausible ("innocent") or implausible ("guilty")
- If the sample mean from our data is very different from the null mean, then we conclude that the null hypothesis is implausible
- How big a difference do we need to observe to refute the null?

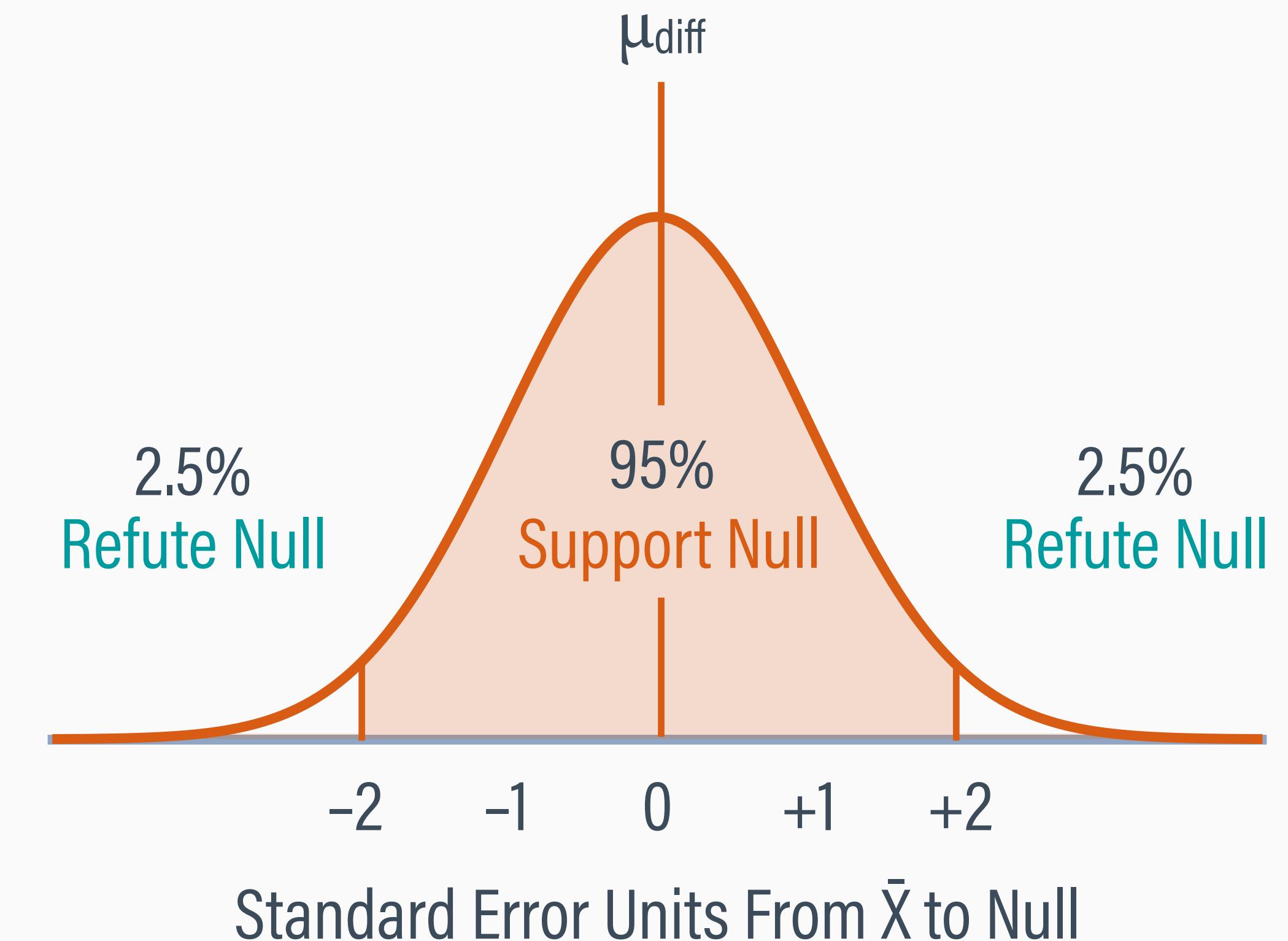
EVALUATING THE NULL

- Any \bar{X}_{diff} near the middle of the sampling distribution ($\mu_{\text{diff}} = 0$) lends support to the null
- Such a sample has a high probability of originating from the null population
- We refute the null if the sample \bar{X}_{diff} falls far from μ_{diff}
- Such a sample has a low probability of originating from the null population



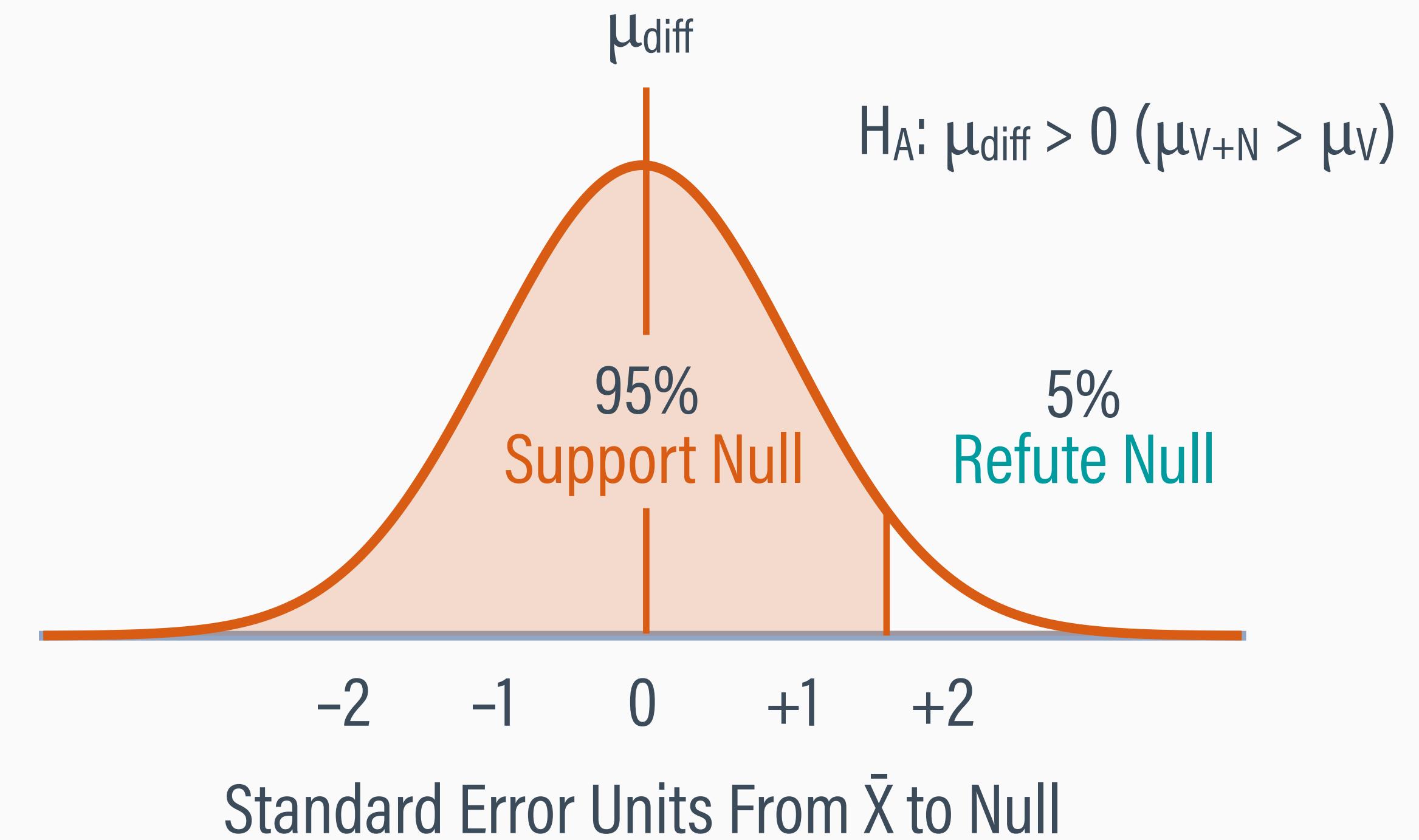
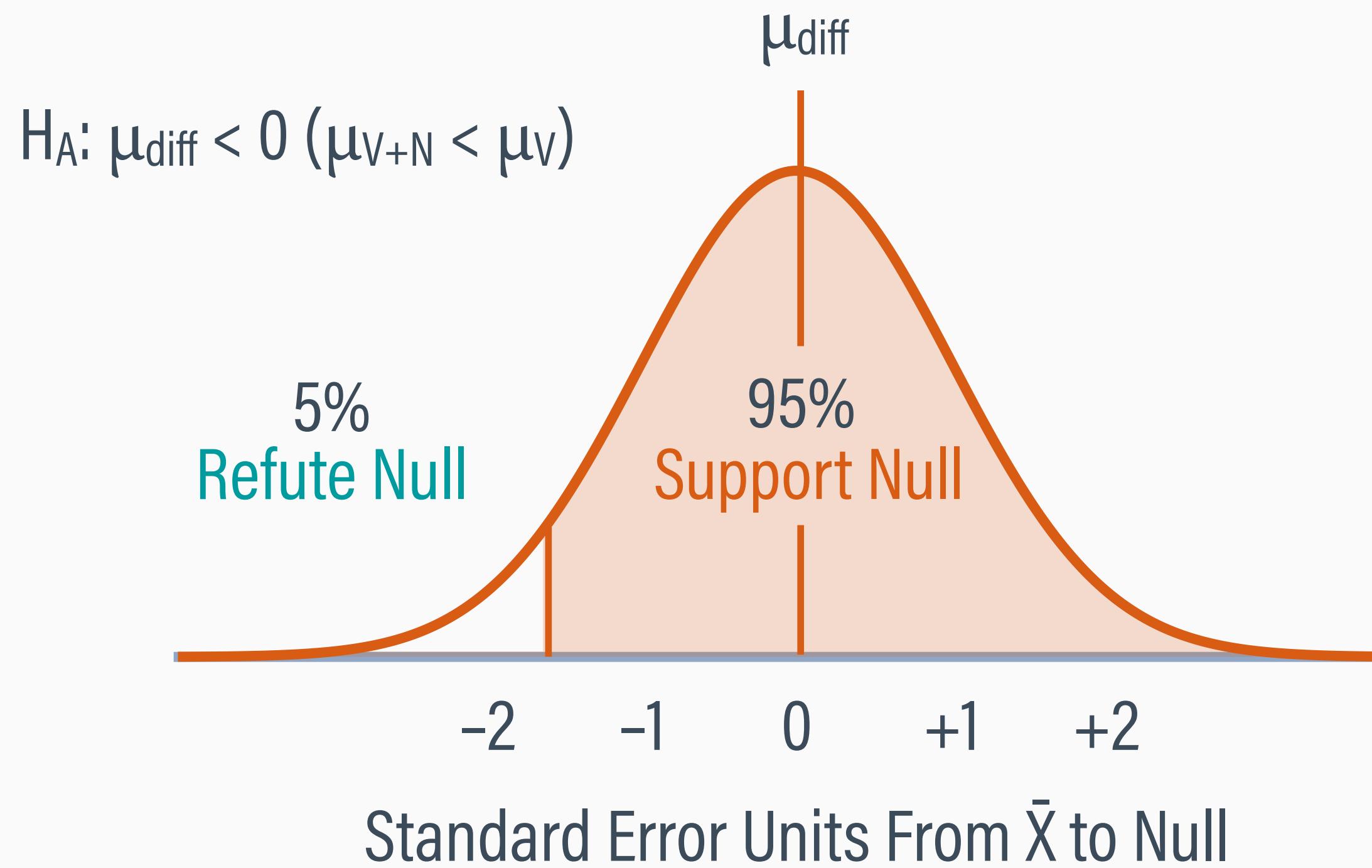
TWO-TAILED ALTERNATE HYPOTHESES

- By convention, we refute the null if the sample \bar{X}_{diff} falls outside the middle 95% of the sampling distribution
- Such a sample has less than a 5% chance of originating from the null population ($p < .05$)
- The 5% rejection region (**alpha level**) is split in half to allow for the possibility that either an increase or a decrease provides evidence against H_0



ONE-TAILED ALTERNATE HYPOTHESES

- The 5% rejection region (**alpha level**) is placed in one tail, since only a positive (or negative) \bar{X}_{diff} counts as evidence against H_0

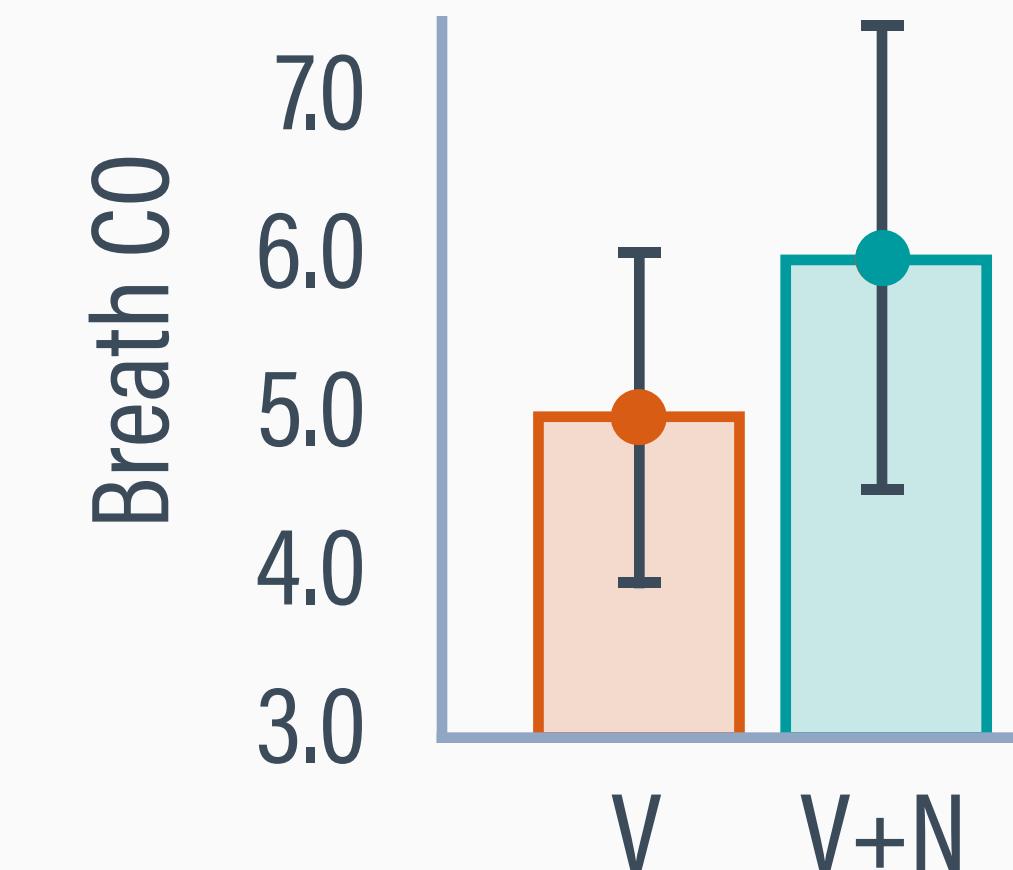


SIGNIFICANCE TESTING STEPS

- 1 Specify hypotheses
- 2 Define standard of evidence
- 3 Design study and collect data
- 4 Compare data to null hypothesis
- 5 Evaluate hypotheses and draw conclusion

ANALYSIS SUMMARY

- $N = 165$ participants across two conditions
- The sample mean difference is $\bar{X}_V - \bar{X}_{V+N} = -1.00$ breath CO units ($5.02 - 6.02 = -1.00$)
- Could a mean difference of 1 point originate from a null population where both groups are equal?



Medication	\bar{X}	s	n
Varenicline	5.02	4.88	82
Varenicline + Naltrexone	6.02	6.86	83

R OUTPUT

Descriptive statistics by group

Condition: **Varenicline**

	vars	n	mean	sd	median	trimmed	mad	min	max	range	skew	kurtosis	se
COWeek8	1	82	5.02	4.88	3.5	4.33	3.71	0	26	26	1.57	3.09	0.54

Condition: **Varenicline + Naltrexone**

	vars	n	mean	sd	median	trimmed	mad	min	max	range	skew	kurtosis	se
COWeek8	1	83	6.02	6.86	3	4.81	4.45	0	29	29	1.62	2.31	0.75

EFFECT SIZE EXAMPLE

- The effect size expresses the mean difference ($\bar{X}_{\text{diff}} = -1$) on a standardized metric
- The mean difference of 1.0 on the breath CO scale equates to 0.17 standard deviation units
- The two means differ by 0.17 z-score units

Medication	\bar{X}	s	n
Varenicline	5.02	4.88	82
Varenicline + Naltrexone	6.02	6.86	83

$$d = \frac{|\bar{X}_1 - \bar{X}_2|}{\sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}}} =$$

$$\frac{|5.02 - 6.02|}{\sqrt{\frac{(82 - 1)4.88^2 + (83 - 1)6.86^2}{82 + 83 - 2}}} = 0.17$$

R OUTPUT

Cohen d statistic of difference between two means

lower effect upper

COWeek8 -0.14 **0.17** 0.47

Multivariate (Mahalanobis) distance between groups

[1] 0.17

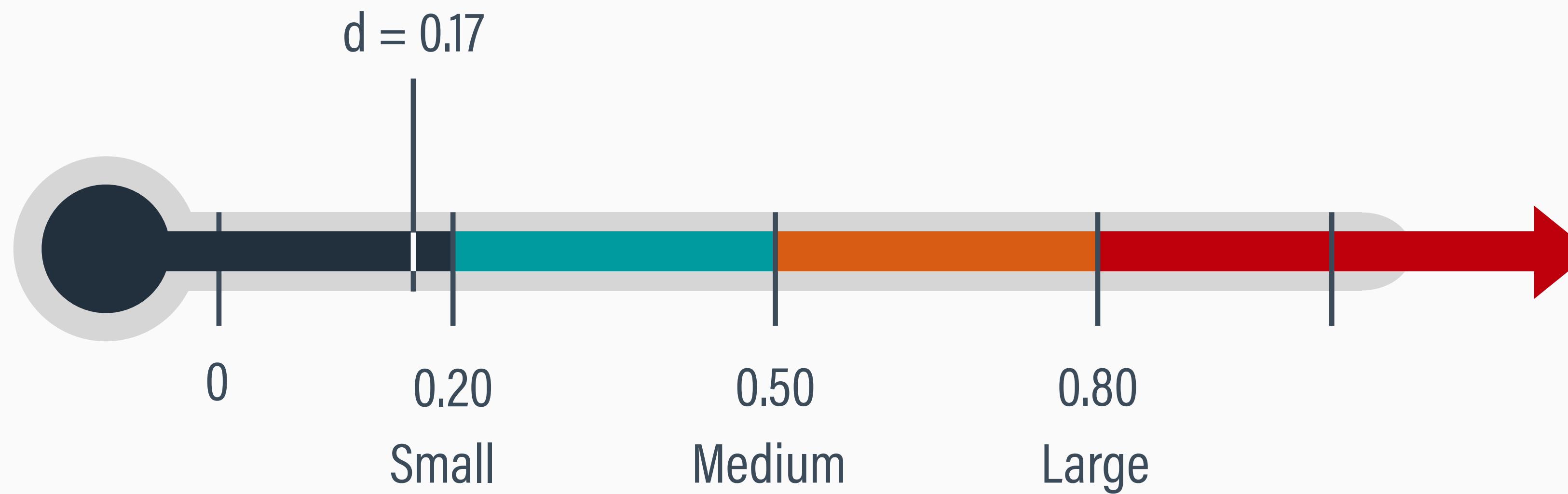
r equivalent of difference between two means

COWeek8

0.08

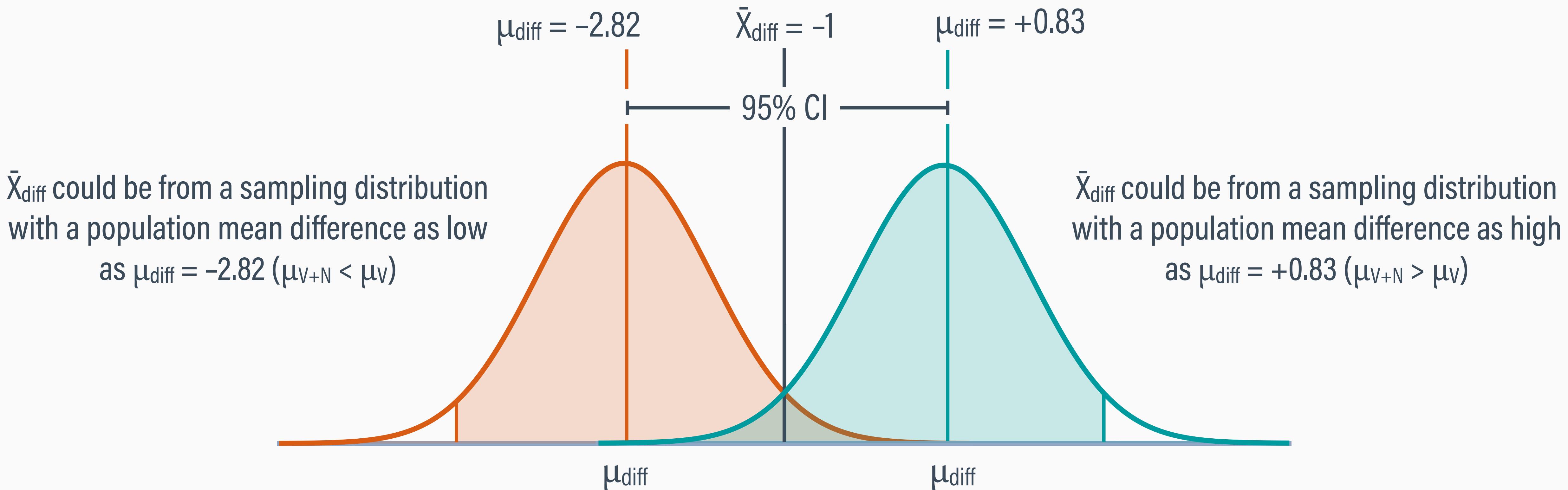
COHEN'S BENCHMARKS

- A 0.17 standard deviation difference is very close to the small effect size cutoff



95% CONFIDENCE INTERVAL

- The 95% confidence interval gives the two most extreme values of the population mean that could have reasonably produced these data



R OUTPUT

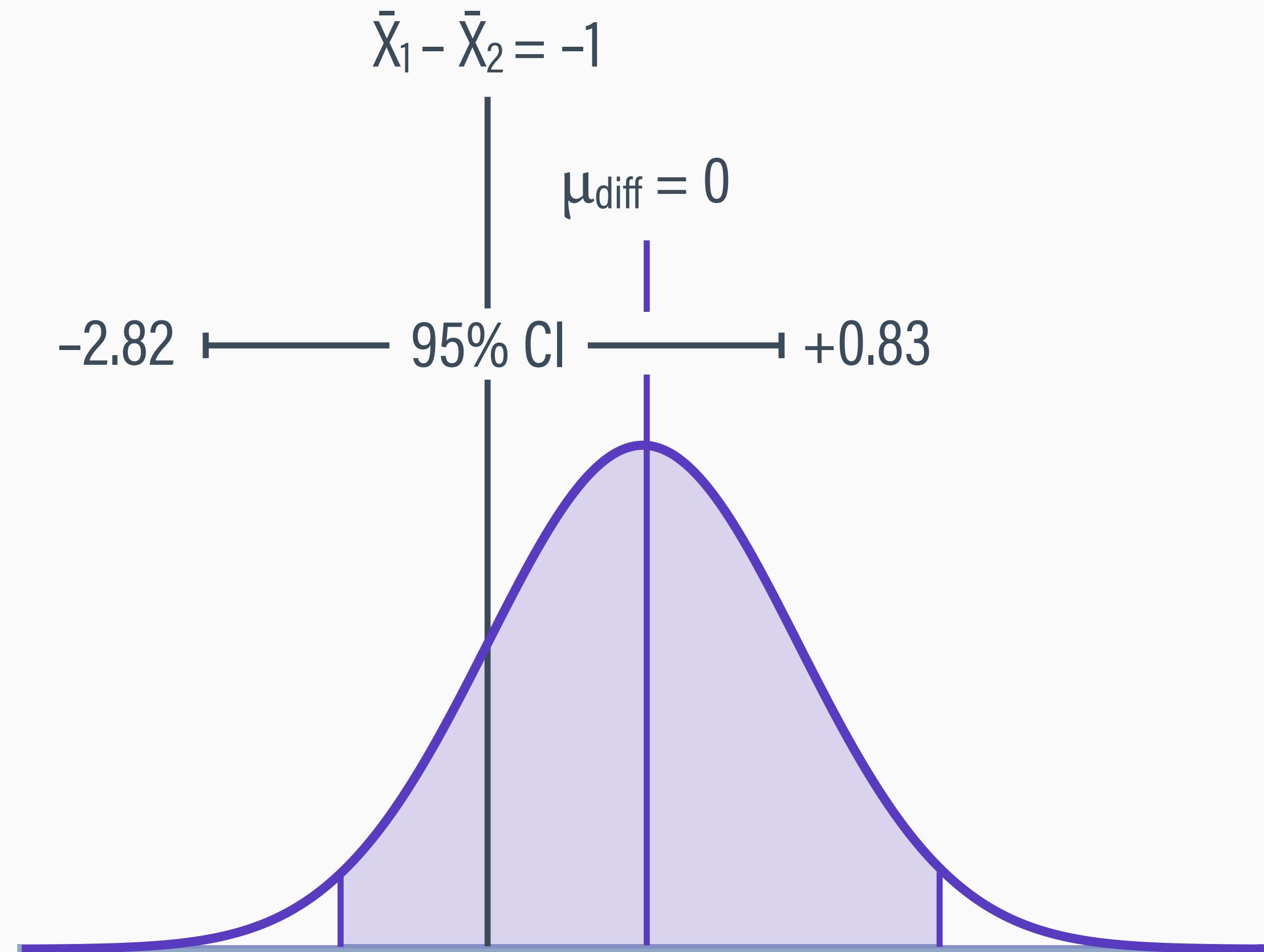
```
data: COWeek8 by Condition  
t = -1.0802, df = 148.19, p-value = 0.2818  
alternative hypothesis: true difference in means between group Varenicline and  
group Varenicline + Naltrexone is not equal to 0  
  
95 percent confidence interval:  
-2.8285605 0.8291482  
  
sample estimates:  
mean in group Varenicline mean in group Varenicline + Naltrexone  
5.024390 6.024096  
  
standard error of mean difference: 0.9254862
```



The study produced a mean difference and 95% confidence interval of $\bar{X}_{\text{diff}} = -1$ and $\text{CI}_{95\%} = [-2.82, +0.83]$. In small groups of two or three, discuss whether this sample of $N = 165$ participants could have reasonably originated from a population where there is truly no mean difference between medication arms ($\mu_{\text{diff}} = 0$).

SIGNIFICANCE TESTING WITH 95% INTERVALS

- A population with a mean difference of 0 is likely to have produced this sample because the null mean is within the 95% interval
- The 95% confidence interval provides the same conclusion as a two-tailed significance test with a .05 significance criterion!



SIGNIFICANCE TESTING STEPS

- 1 Specify hypotheses
- 2 Define standard of evidence
- 3 Design study and collect data
- 4 Compare data to null hypothesis
- 5 Evaluate hypotheses and draw conclusion

COMPARING DATA TO THE NULL

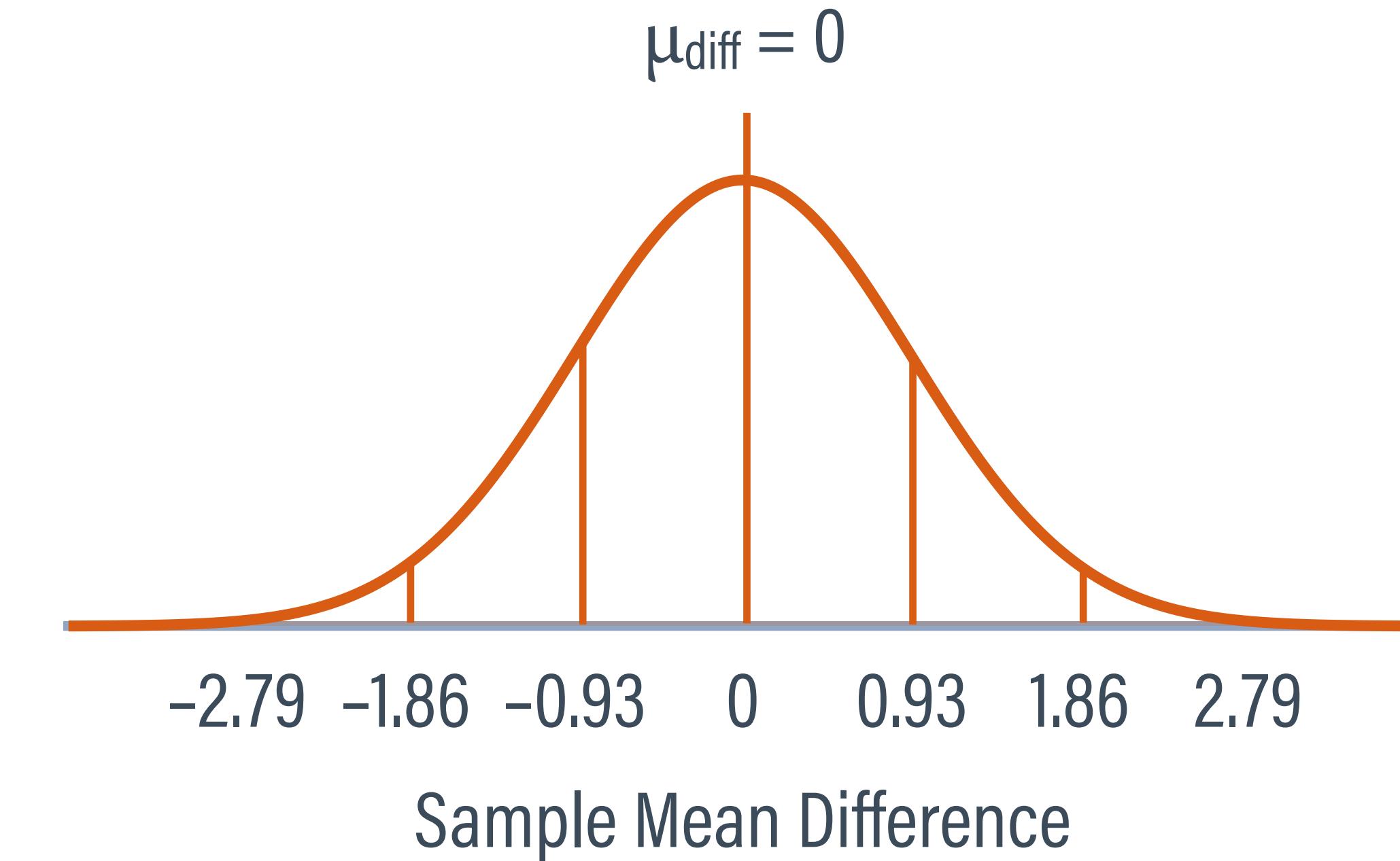
- Two ways to determine whether the sample \bar{X}_{diff} is consistent (or inconsistent) with the null population mean
- The t-statistic gives a standardized distance between the sample mean and the null hypothesis mean (like a z-score)
- A p-value tells us how likely it is that hypothetical samples like our data would originate from the null population

STANDARD ERROR OF A MEAN DIFFERENCE

- The standard error gives the expected (average) sampling error in a mean difference statistic across many hypothetical samples

$$s_{\bar{x}_{\text{diff}}} = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} = \sqrt{\frac{4.88^2}{82} + \frac{6.86^2}{83}} = 0.93$$

- Across many hypothetical samples from a null population, we would expect the mean difference to be ± 0.93 breath CO points from $\mu_{\text{diff}} = 0$



Consider the sampling distribution of sample means from a null population with $\mu_{\text{diff}} = 0$. The sample mean difference of $\bar{X}_{\text{diff}} = -1$ ($s_{\bar{X}_{\text{diff}}} = 0.93$). In small groups of two or three, discuss whether the data provide evidence for or against the null hypothesis.

t-STATISTIC

- The t-statistic quantifies the number of standard error units that separate the sample mean and null hypothesis population mean

$$t = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{s_{\bar{X}_{\text{diff}}}} = \frac{\text{distance from the null}}{\text{standard error (std. dev. of } \bar{X}_{\text{diff}})}$$

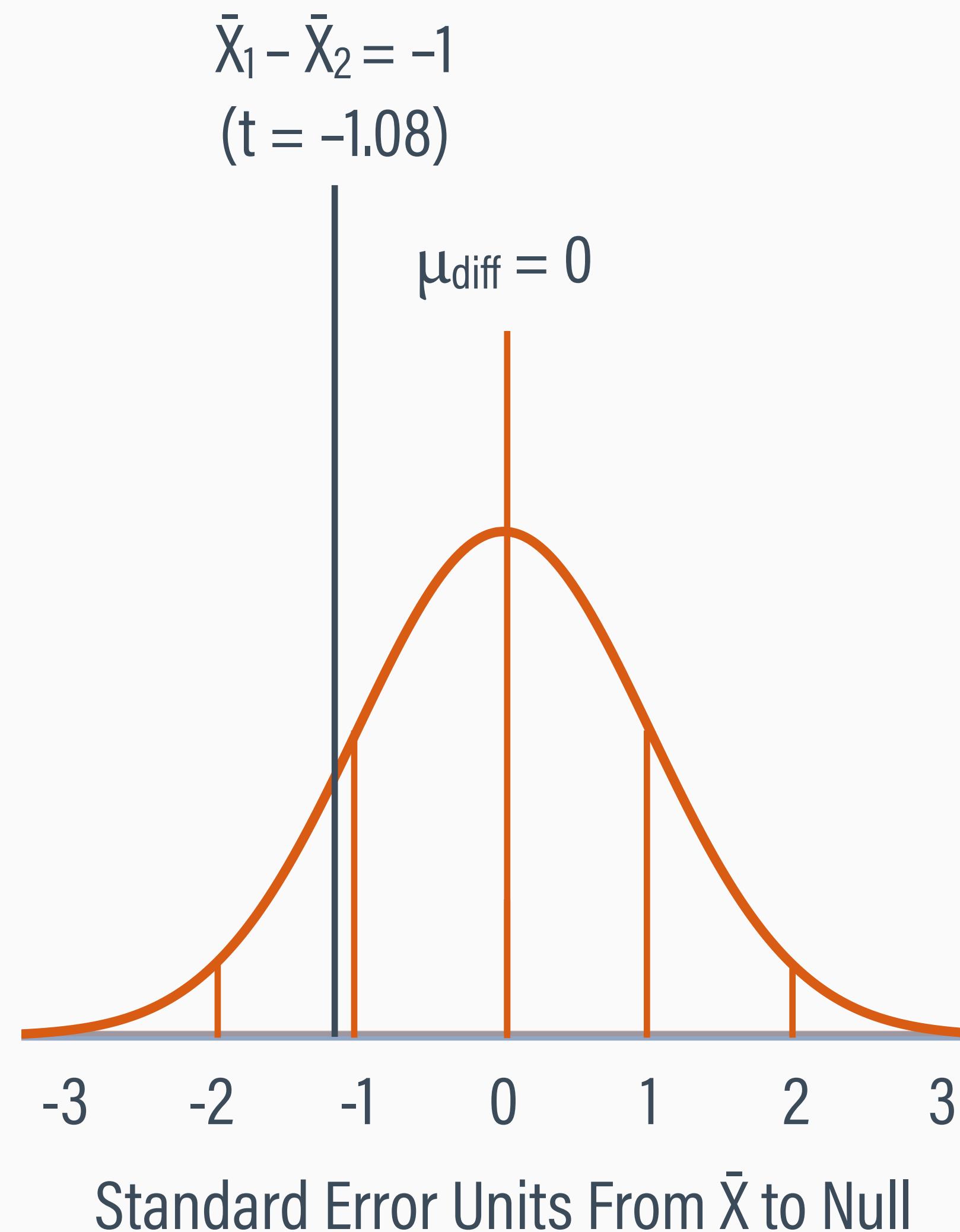
- The t-statistic is the same as a z-score (a standardized metric where distance is expressed in standard deviation units)

t-STATISTIC EXAMPLE

- The t-statistic indicates that 1.08 standard error units separate the sample mean and null

$$t = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{s_{\bar{X}_{\text{diff}}}} = \frac{(5.02 - 6.02) - 0}{0.93} = -1.08$$

- The negative sign is a result of subtracting the higher mean from the lower mean

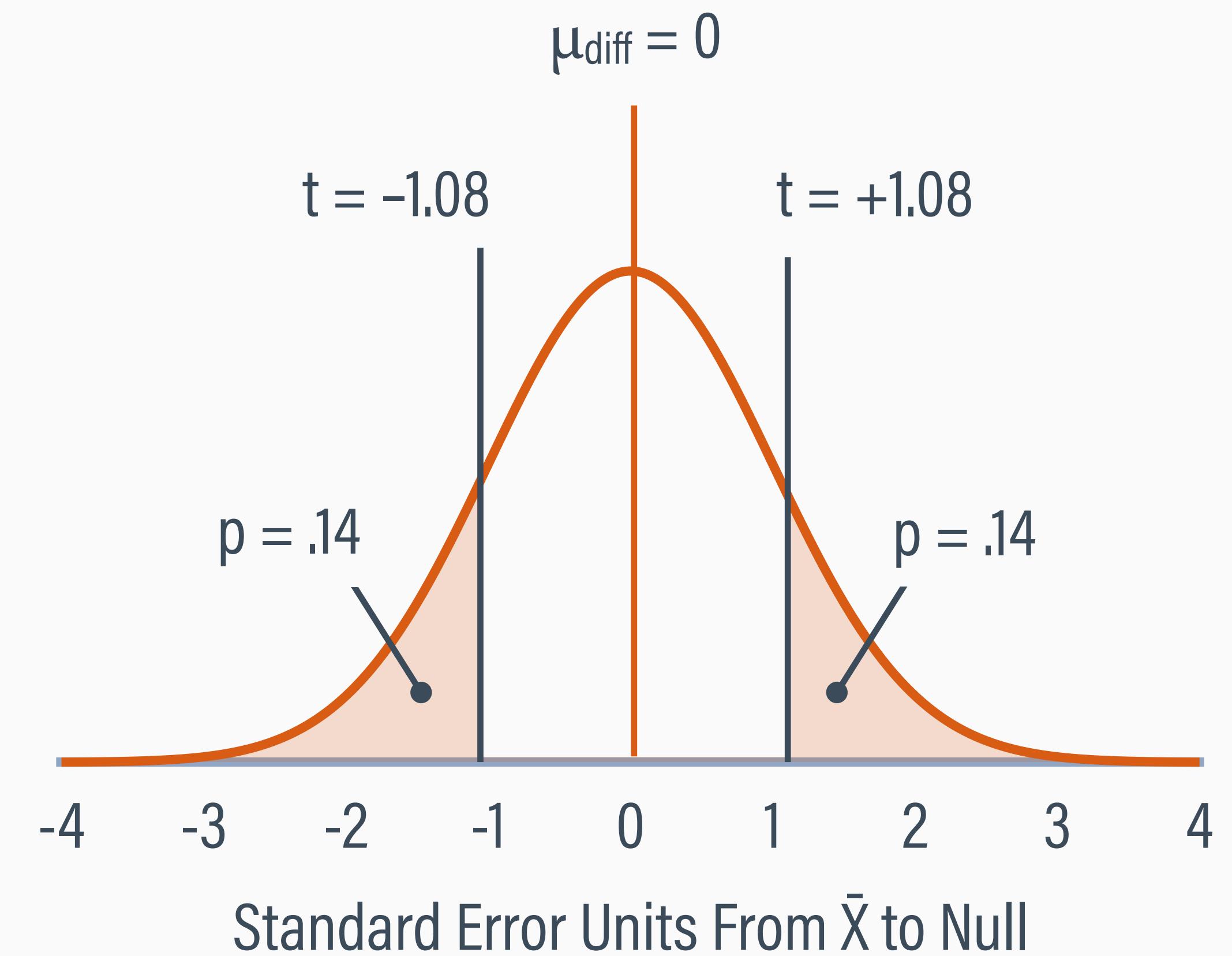


PROBABILITY VALUES (P-VALUES)

- A p-value is defined as proportion of hypothetical samples that have a t-statistic at least as large as the sample data
- Assuming the null is true, how likely is it to draw a sample with an effect at least as large as the one from our data?
- Visually, probability is an area under the curve, obtained by applying calculus integrals to the t-distribution function

TWO-TAILED P-VALUE

- The p-value tells how likely it is to draw a sample mean difference at least as extreme as ours from a null population with $\mu_{\text{diff}} = 0$
- The probability of drawing a sample from the null population with a t-statistic of at least ± 1.08 is $p = .28$
- 28% of all hypothetical samples from a null population would have t-statistics this large





Suppose the researchers had instead specified a one-tailed test where the predicted that the Varenicline-only condition would have a lower mean (i.e., they correctly predicted the direction). In small groups of two or three, discuss how the p-value would change with a one-tailed alternate hypothesis.

R OUTPUT

data: COWeek8 by Condition

t = -1.0802, df = 148.19, p-value = 0.2818

alternative hypothesis: true difference in means between group Varenicline and group Varenicline + Naltrexone is not equal to 0

95 percent confidence interval:

-2.8285605 0.8291482

sample estimates:

mean in group Varenicline	mean in group Varenicline + Naltrexone
5.024390	6.024096

standard error of mean difference: 0.9254862

SIGNIFICANCE TESTING STEPS

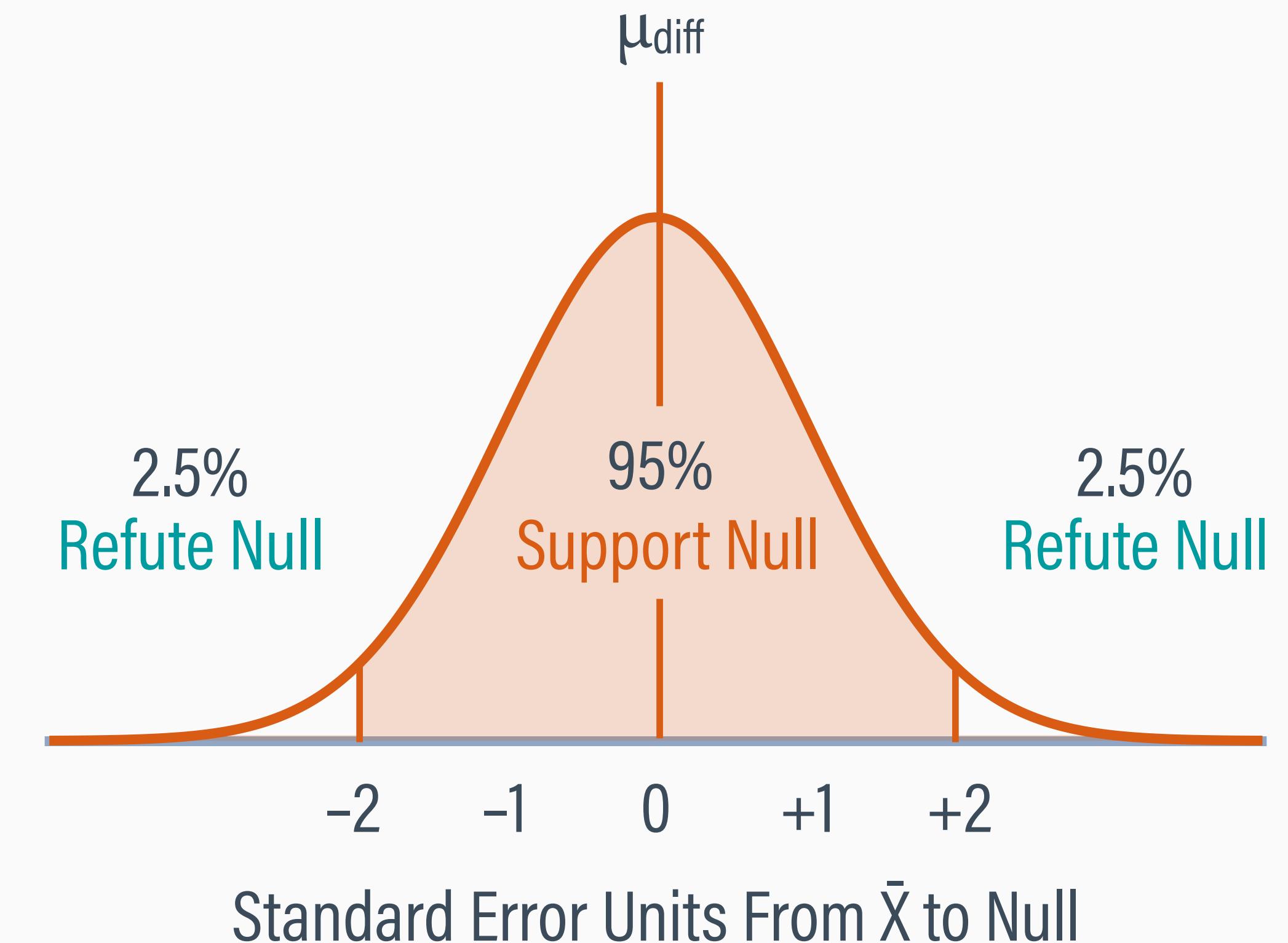
- 1 Specify hypotheses
- 2 Define standard of evidence
- 3 Design study and collect data
- 4 Compare data to null hypothesis
- 5 Evaluate hypotheses and draw conclusion

RESEARCH QUESTION REVISITED

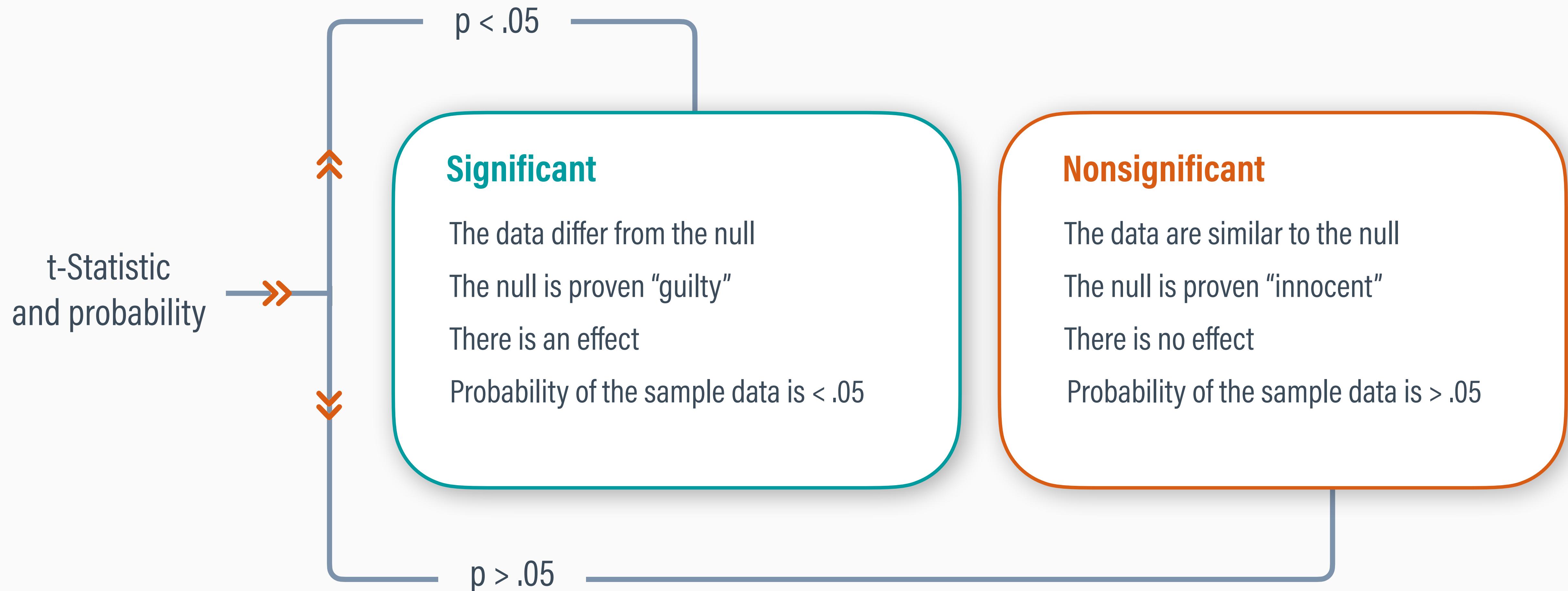
- Studies typically attempt to answer a handful of research questions involving associations between key variables
- Does smoking intensity differ between participants receiving a combination of two medications and those receiving a single medication alone?
- The null (no effect) hypothesis states that the breath CO means are identical (the population mean difference is zero)

5% SIGNIFICANCE CRITERION REVISITED

- By convention, we refute the null if the sample \bar{X}_{diff} falls outside the middle 95% of the sampling distribution ($p < .05$)
- Such a sample has less than a 5% chance of originating from the null population
- We deem the null implausible because our data are unlikely to originate from that population

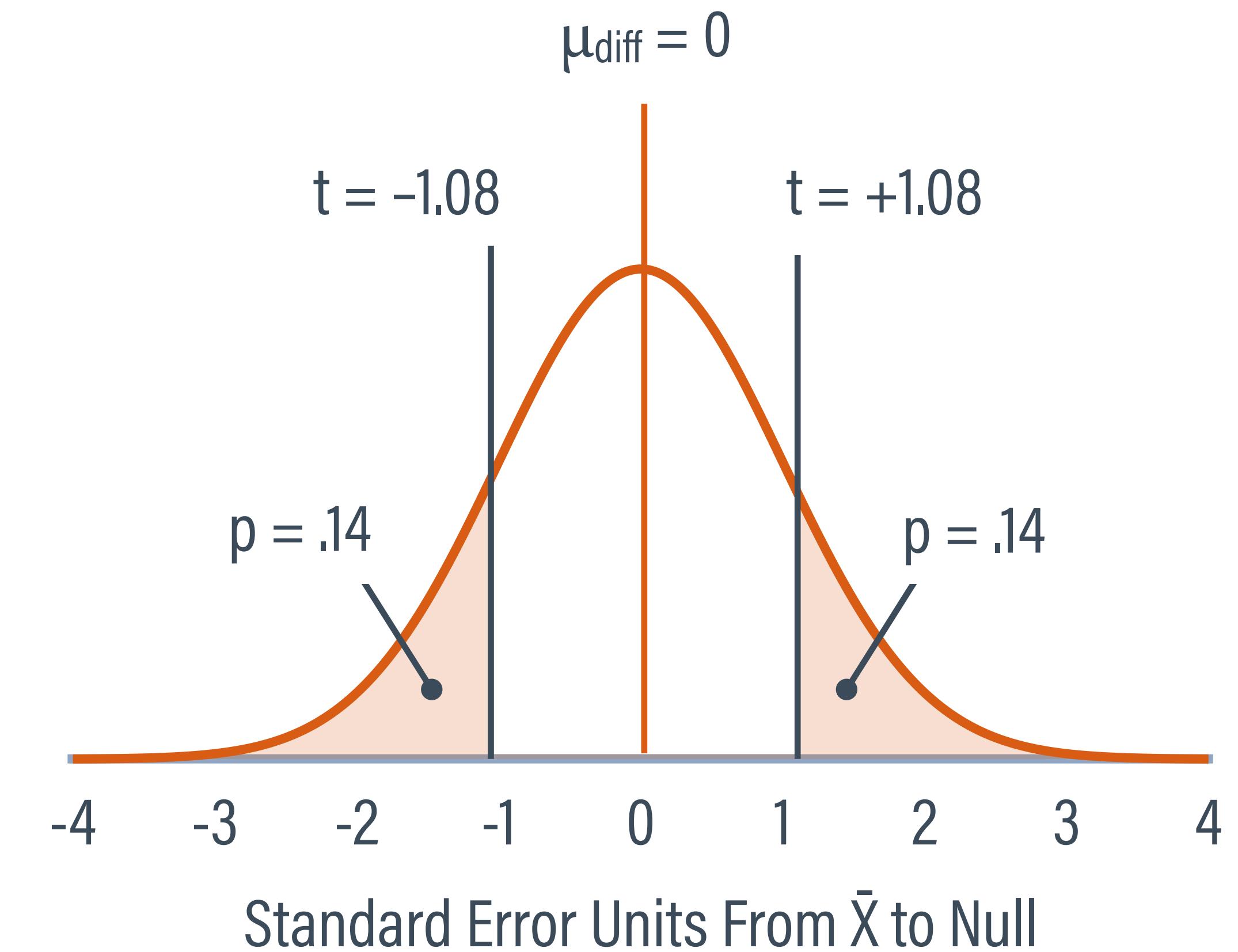


DECISION TREE





The two-tailed probability for the study is $p = .28$. In small groups of two or three, discuss your decision about the null hypothesis. Translate your decision into a tangible statement about the effect of the two medications on smoking intensity.



CONCLUSION: TWO-TAILED ALTERNATE

- The p-value of .28 would lead us to retain (support) the null
- A mean difference as large as $\bar{X}_{\text{diff}} = \pm 1$ breath CO points (or equivalently, a t-statistic at least ± 1.08) is quite likely to have originated from a null population with $\mu_{\text{diff}} = 0$
- The sample mean difference provides evidence that the two medication arms could be equivalent in the population

FALSE NEGATIVE (TYPE II ERRORS)

- There is always a possibility of committing an inferential error (making the wrong decision about the null)
- Even when \bar{X}_{diff} falls within the middle 95% of the sampling distribution, it could still come from a population with a true difference (the effect may be too small to detect with our N)
- We conclude there is no difference, while acknowledging the possibility of a false negative—supporting the null when it is actually false (a Type II error)

APA-STYLE ANALYSIS SUMMARY

We used a independent-samples *t*-test to examine whether the two medication regimes produced different impacts on smoking cessation, as measured by breath CO. Table 1 gives the descriptive statistics. The mean difference was approximately one point (single medication lower), with a 95% confidence interval for the mean difference that ranged from -2.83 (single medication superior) to 0.83 (dual medication superior). An independent t-test revealed a non-significant difference between the two medication arms, $t(163) = -1.08, p = .28$. Finally, the standardized mean difference was just below Cohen's small effect size benchmark ($d = 0.17$), indicating a subtle mean difference in the sample data.

OUTLINE

- 1 Between-group designs
- 2 Quick review
- 3 Significance testing steps
- 4 Statistical assumptions
- 5 Study questions

STATISTICAL ASSUMPTIONS

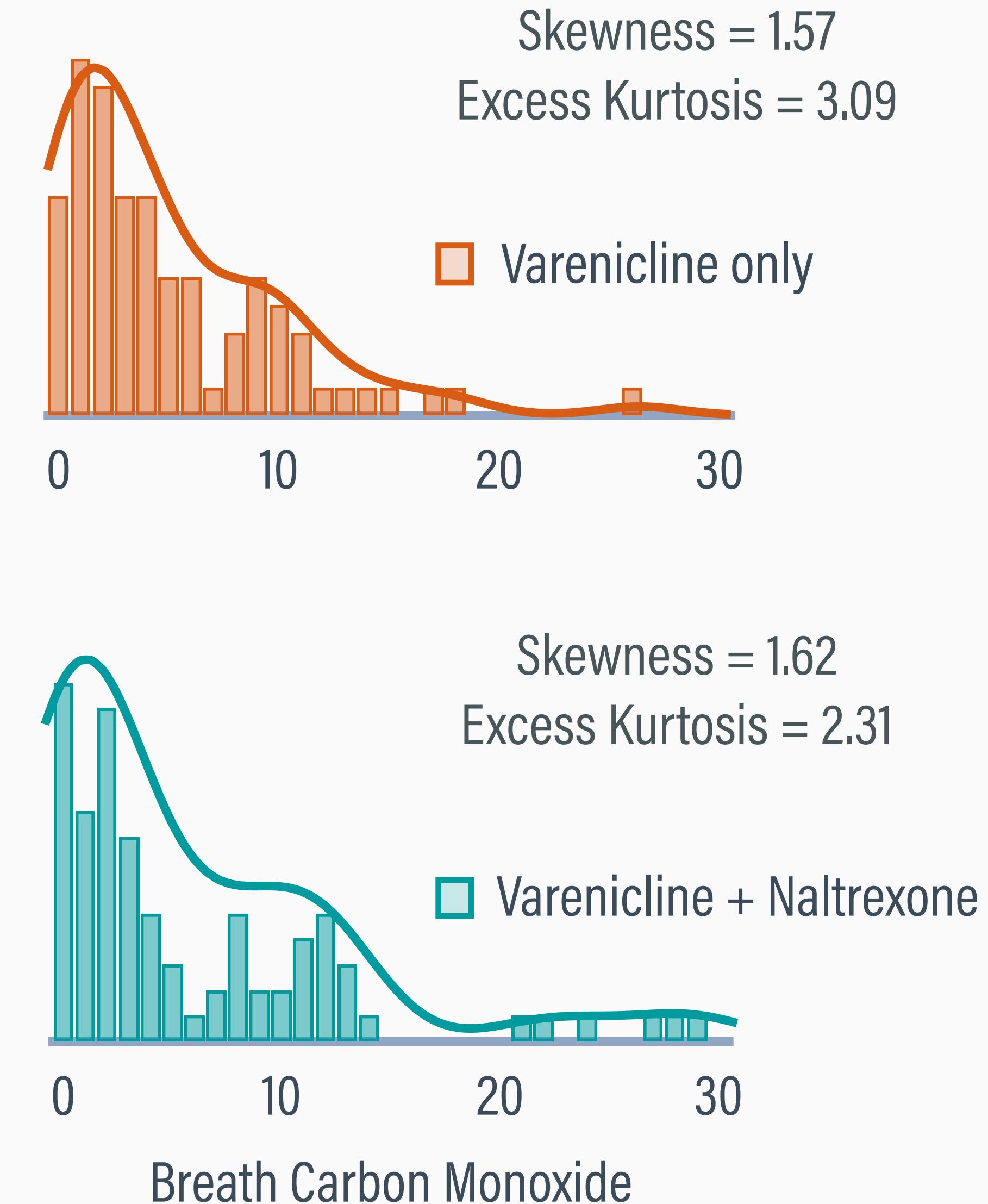
- The accuracy of t-tests (and other statistics) depends on certain conditions in the data being true (e.g., normality)
- Violations of assumptions can bias estimates, inflate or deflate standard errors, and distort significance tests
- Always check reasonableness of assumptions before drawing conclusions

INDEPENDENT T-TEST ASSUMPTIONS

- Numeric (approximately continuous) dependent variable
- Scores are approximately normal in each population
- Independence of observations (no participant's score influences any other participant's score)
- The two populations have equal variances (homogeneity of variance)

WITHIN-GROUP NORMALITY

- In small samples, normality violations can artificially inflate or deflate standard errors, thus distorting significance tests
- Normality is less of a concern if the sample size is large enough (e.g., $N_s > 40$ to 50)
- Normalizing transformations (e.g., the natural log of scores) are common in some domains



HOMOGENEITY OF VARIANCE

- The sample variances (squared standard deviations) differ by a factor of about two ($4.88^2 = 23.81$ versus 6.86^2)
- The classic **Student's t-test** assumes equal variation
- The **Welch t-test** we use relaxes this assumption, adjusting both the standard error and degrees of freedom
- Methodological literature shows Welch's test is quite robust to unequal variances and sample sizes

Medication	\bar{x}	s	n
Varenicline	5.02	4.88	82
Varenicline + Naltrexone	6.02	6.86	83

OUTLINE

- 1 Between-group designs
- 2 Quick review
- 3 Significance testing steps
- 4 Statistical assumptions
- 5 Study questions

STUDY QUESTIONS

A researcher wants to determine whether male and female Latinx youth experience different levels of discrimination. To do so, he recruits a sample of $n = 339$ females and $n = 300$ males. The sample mean difference is $\bar{X}_{\text{diff}} = \bar{X}_{\text{female}} - \bar{X}_{\text{male}} = 13.97 - 13.11 = +0.86$. To evaluate whether gender differences exist, you will perform significance testing steps assuming a population with a true mean difference of $\mu_{\text{diff}} = 0$.

STUDY QUESTIONS (1)

1. State the null hypothesis, both as a sentence and using statistical symbols.

2. State the two-tailed alternate hypothesis, both as a sentence and using statistical symbols.

3. Explain why a independent-samples t-test is the appropriate statistical analysis for this scenario.

STUDY QUESTIONS (2)

4. The sampling distribution under the null hypothesis plays a vital role in hypothesis testing with the independent t-test. Explain how the 5% significance criterion is applied to this distribution, and how it is used to decide whether to reject the null hypothesis.

5. The standard error of the mean difference is $s_{\bar{x}_{\text{diff}}} = 0.35$. Explain what the standard error measures. How does it help you gauge whether a mean difference of +0.86 is similar or different from the null?

STUDY QUESTIONS (3)

6. The t-statistic is $t = 2.43$. Explain what the t-statistic measures. What do the sign and the magnitude of the t-statistic indicate about the plausibility of the null hypothesis?
7. Researchers report the results as “statistically significant.” What is your decision about the null hypothesis. Translate your decision into a tangible statement about the difference in male and female discrimination experiences.

STUDY QUESTIONS (4)

8. The two-tailed p-value was .02. Provide an interpretation of the probability value (I am not asking whether the test is significant).

9. The sample mean difference $\bar{X}_{\text{diff}} = \bar{X}_{\text{female}} - \bar{X}_{\text{male}} = 13.97 - 13.11 = +0.86$. The 95% confidence interval limits are $CI_{95\%} = [0.17, 1.55]$. Provide an interpretation of the confidence interval (I am not asking about its statistical properties). Discuss whether the confidence interval supports or refutes the hypothesis that males and females have equal discrimination experiences.