QUANTITATIVE ANALYSIS

# DIFFERENCE OF MEANS (1)

# AGENDA

1. Follow-up

2. Revisiting Distributions

3. One Sample

4. Independent Samples

5. Dependent Samples

6. Effect Sizes

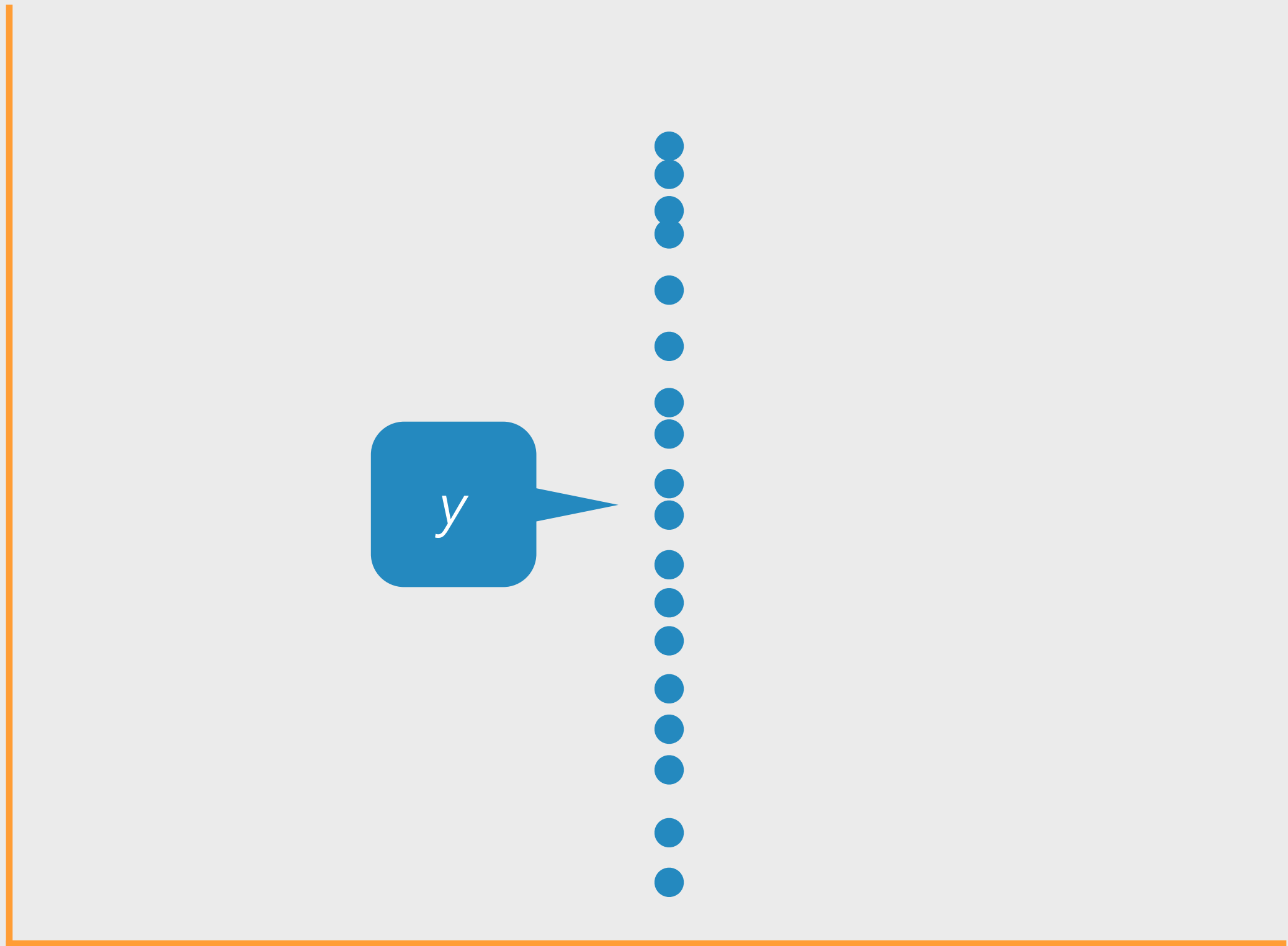# 1 FOLLOW-UP
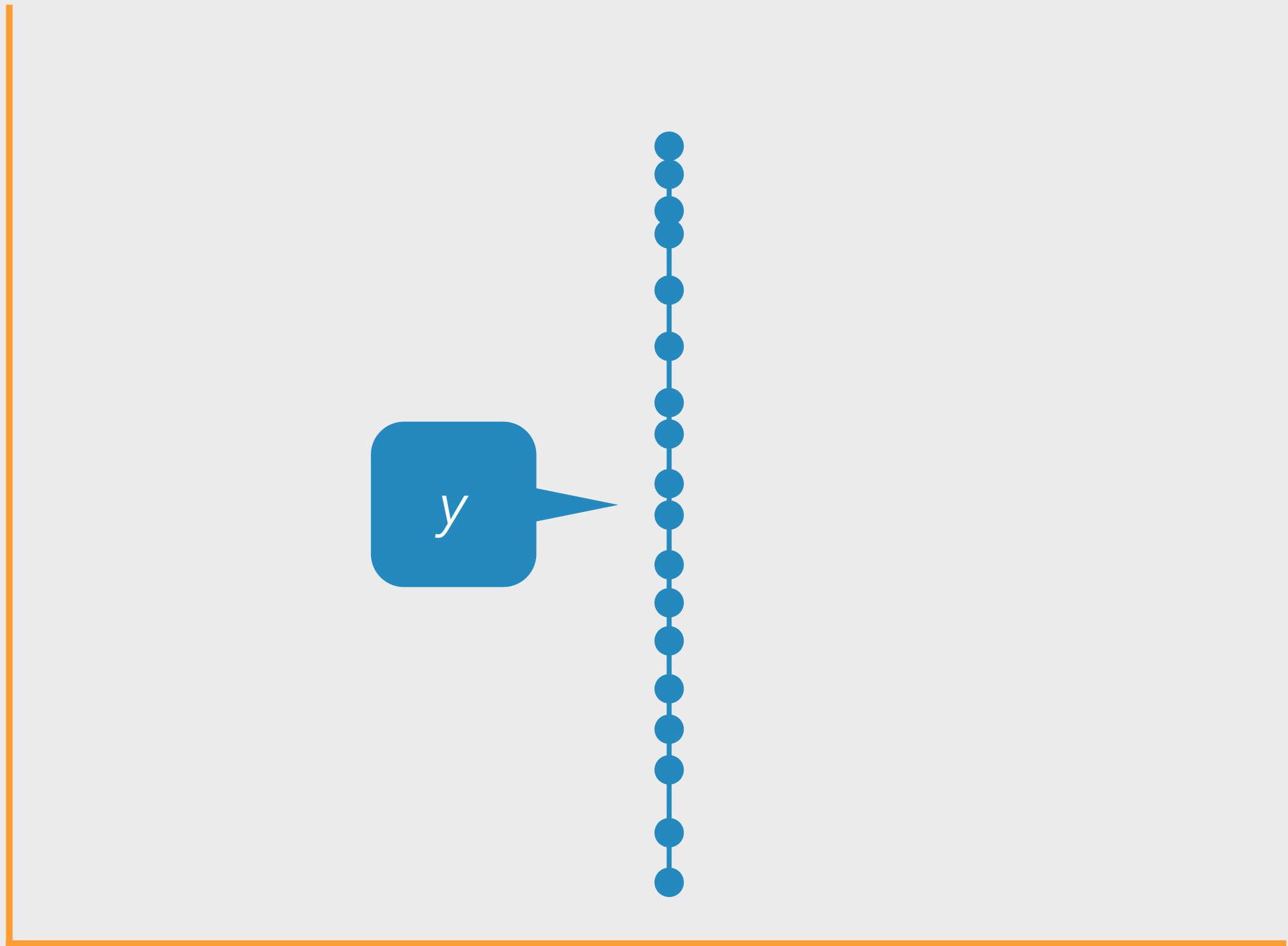
# 2 REVISITING DISTRIBUTIONS

# VARIANCE

**SECOND MOMENT**

**DEFINITION**
SUM OF ALL DEVIANCES, SQUARED AND DIVIDED BY ONE DEGREE OF FREEDOM; EXPECTATION OF HOW DISTRIBUTION DEVIATES FROM THE MEAN
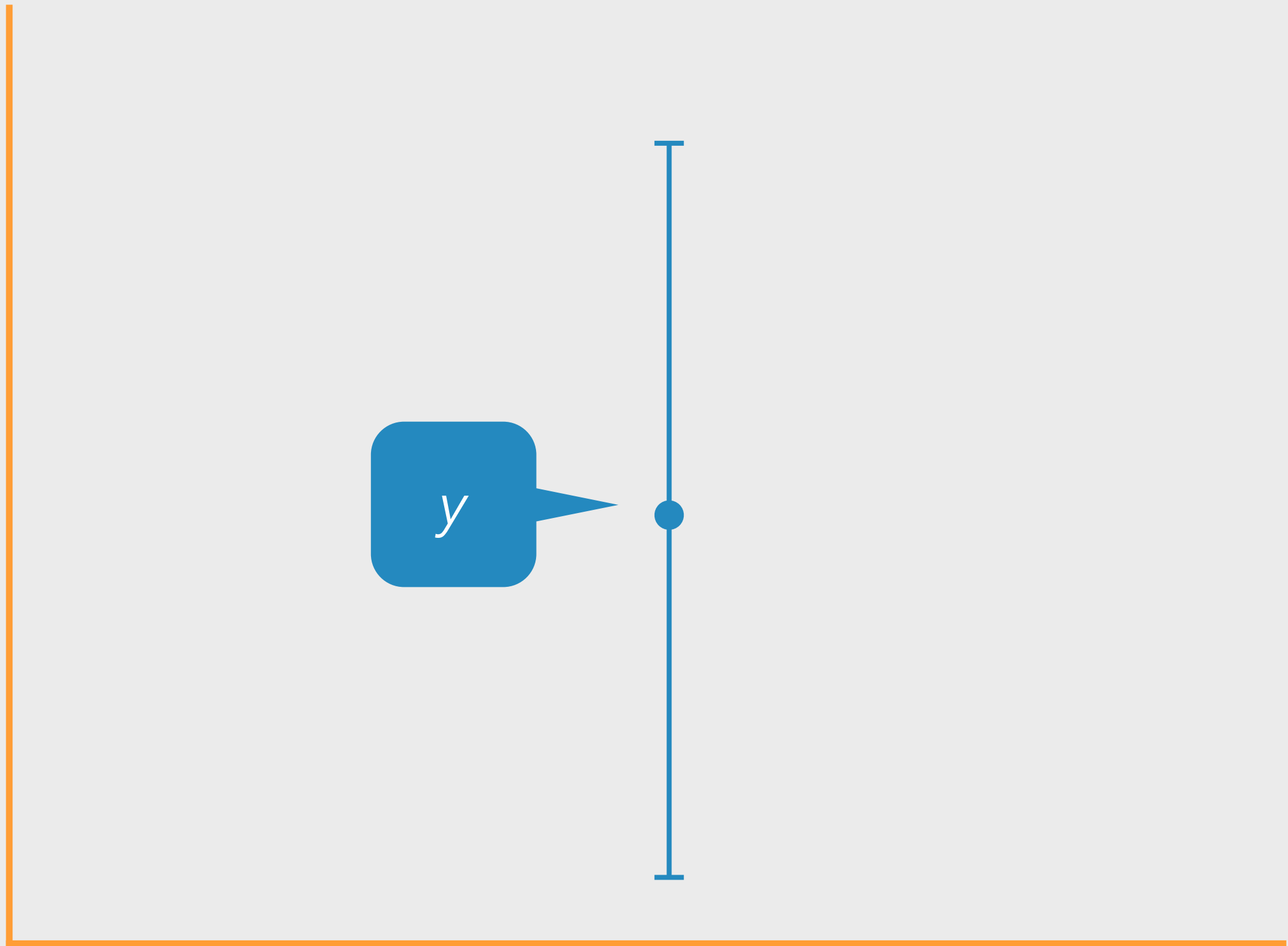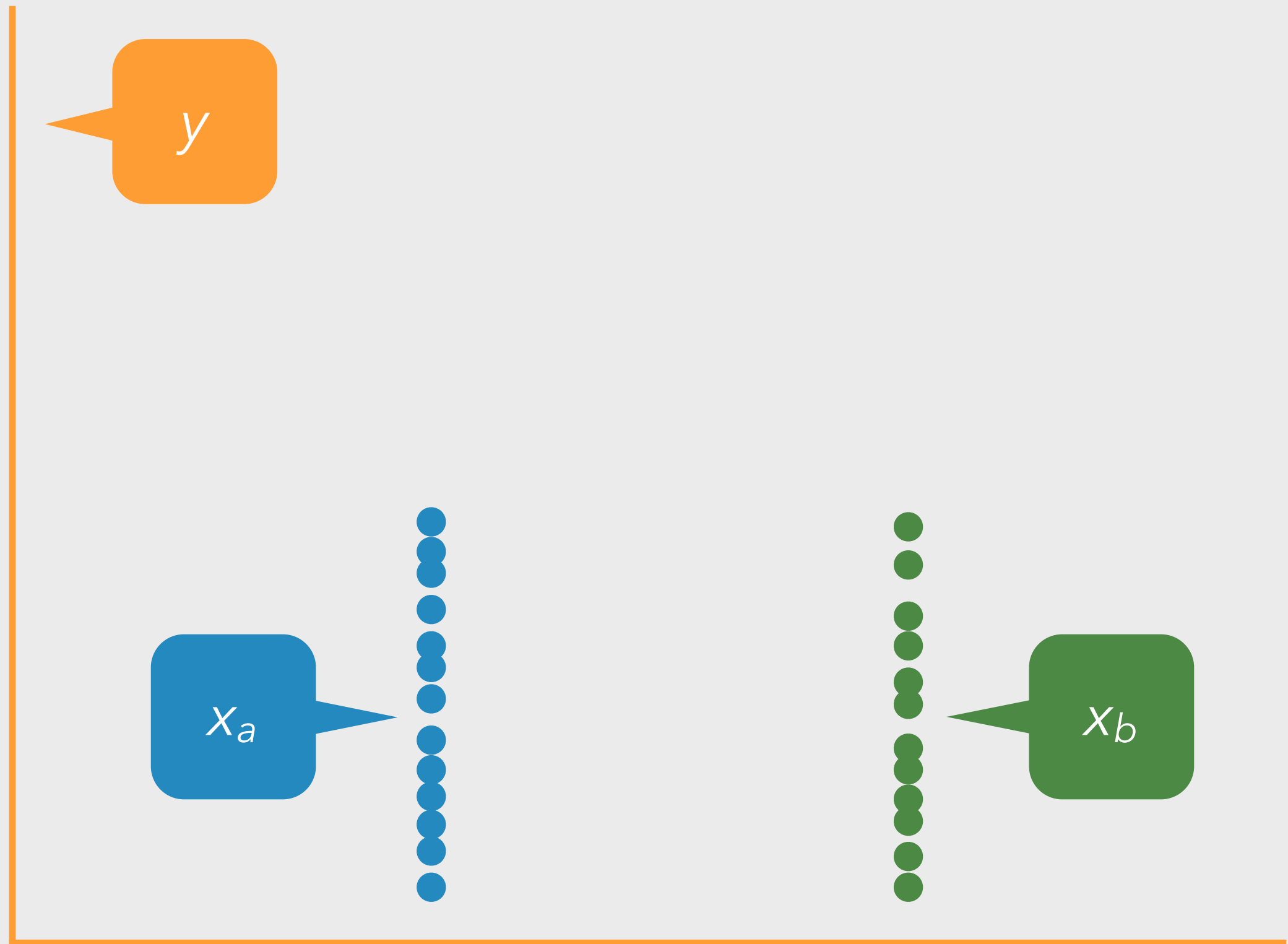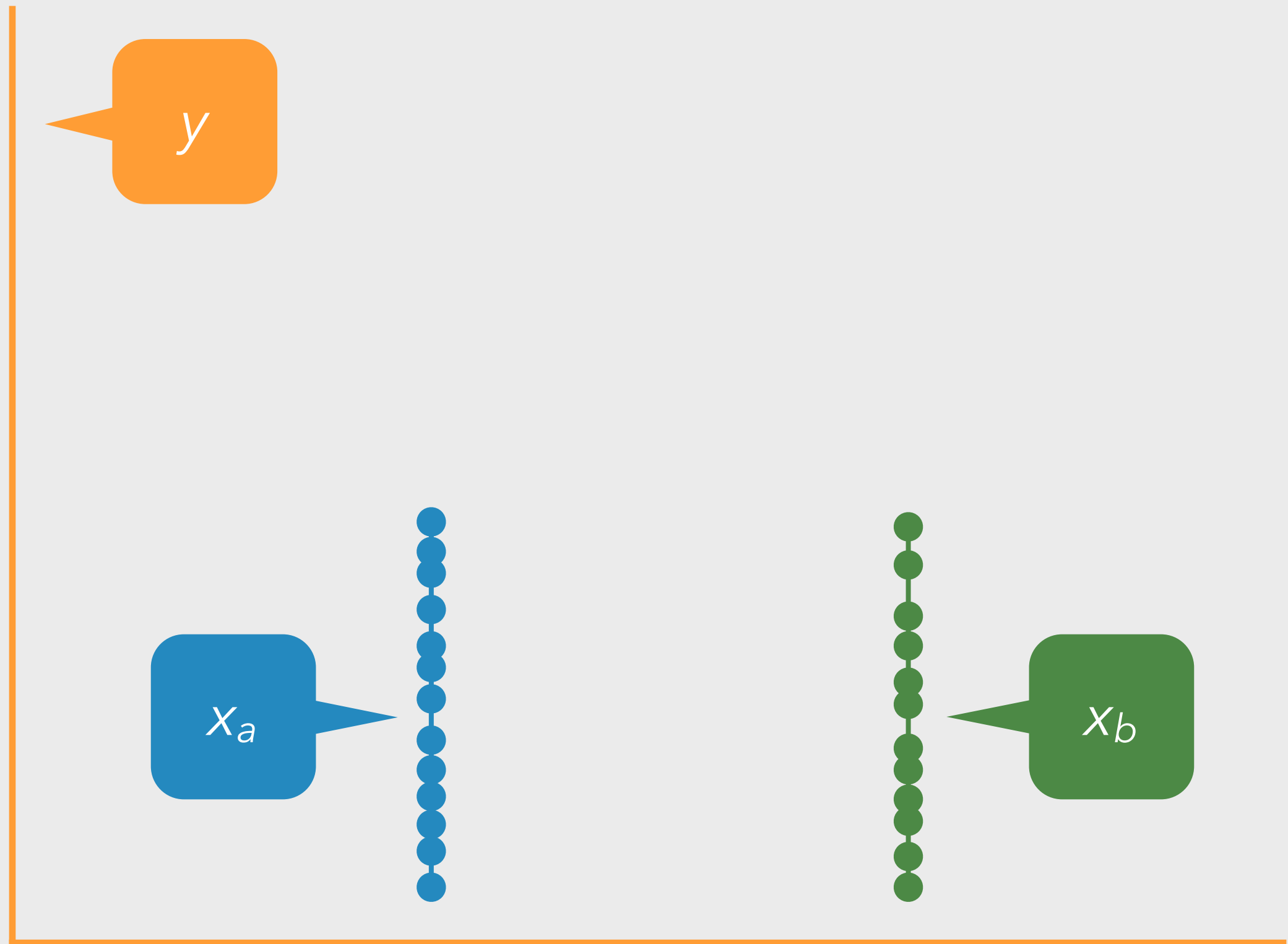
$$s^2 = \frac{\sum_{i=1}^{n}(x - \overline{x})^2}{n - 1}$$

# 3 ONE SAMPLE

# STUDENT'S T-TEST

WILLIAM SEALY GOSSET (1876–1937)
"STUDENT"

▸ Employee of the Guinness company who published his work under the pseudonym "Student".

▸ Student of Karl Pearson's while on research leaves from Guinness.

▸ Original t-tests were developed to conducting quality control testing on Guinness stout.

**Probability Density Functions Compared**
Standard Normal ($\mu=0$, $\sigma=1.0$) and Student's $t$ ($df=1$)

**Probability Density Functions Compared**
Standard Normal ($\mu$=0, $\sigma$=1.0) and Student's $t$ ($df$=20)

Std. Normal)   Student's $t$

**Probability Density Functions Compared**
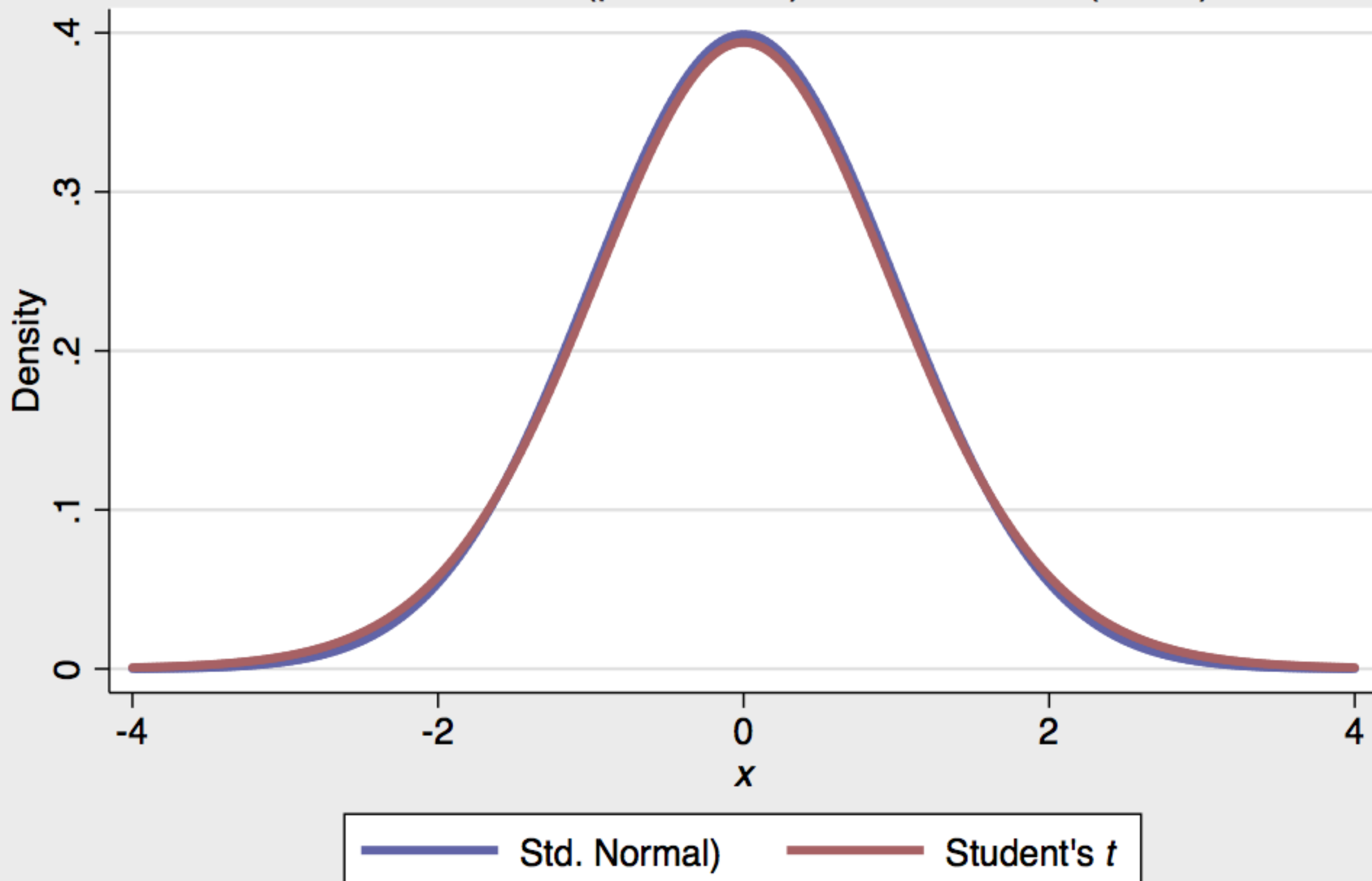Standard Normal ($\mu=0$, $\sigma=1.0$) and Student's $t$ ($df=30$)

Density

$x$

——— Std. Normal)          ——— Student's $t$

# Probability Density Functions Compared
## Standard Normal (μ=0, σ=1.0) and Student's *t* (*df*=100)



Std. Normal)

Student's *t*

# HYPOTHESES

▸ $H_0$ = there is no significant difference between the mean of $y$ and the population


▸ $H_1$ = there is a significant difference between the mean of $y$ and the population

# ASSUMPTIONS

▸ continuous data (*y*)

▸ the distribution of *y* is approximately normal

▸ degrees of freedom (*v*) = *n*-1

# FORMULA

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$$

standard error

# FIND THE PROBABILITY OF T

```
display ttail(df,t)*2


. display ttail(72,3.6308)*2

.0005255


. display ttail(72,1.6308)*2

.1072996
```

# FIND THE PROBABILITY OF T

```
display (1-ttail(df,-t))*2


. display (1-ttail(72,-3.6308))*2

.0005255


. display (1-ttail(72,-1.6308))*2

.1072996
```
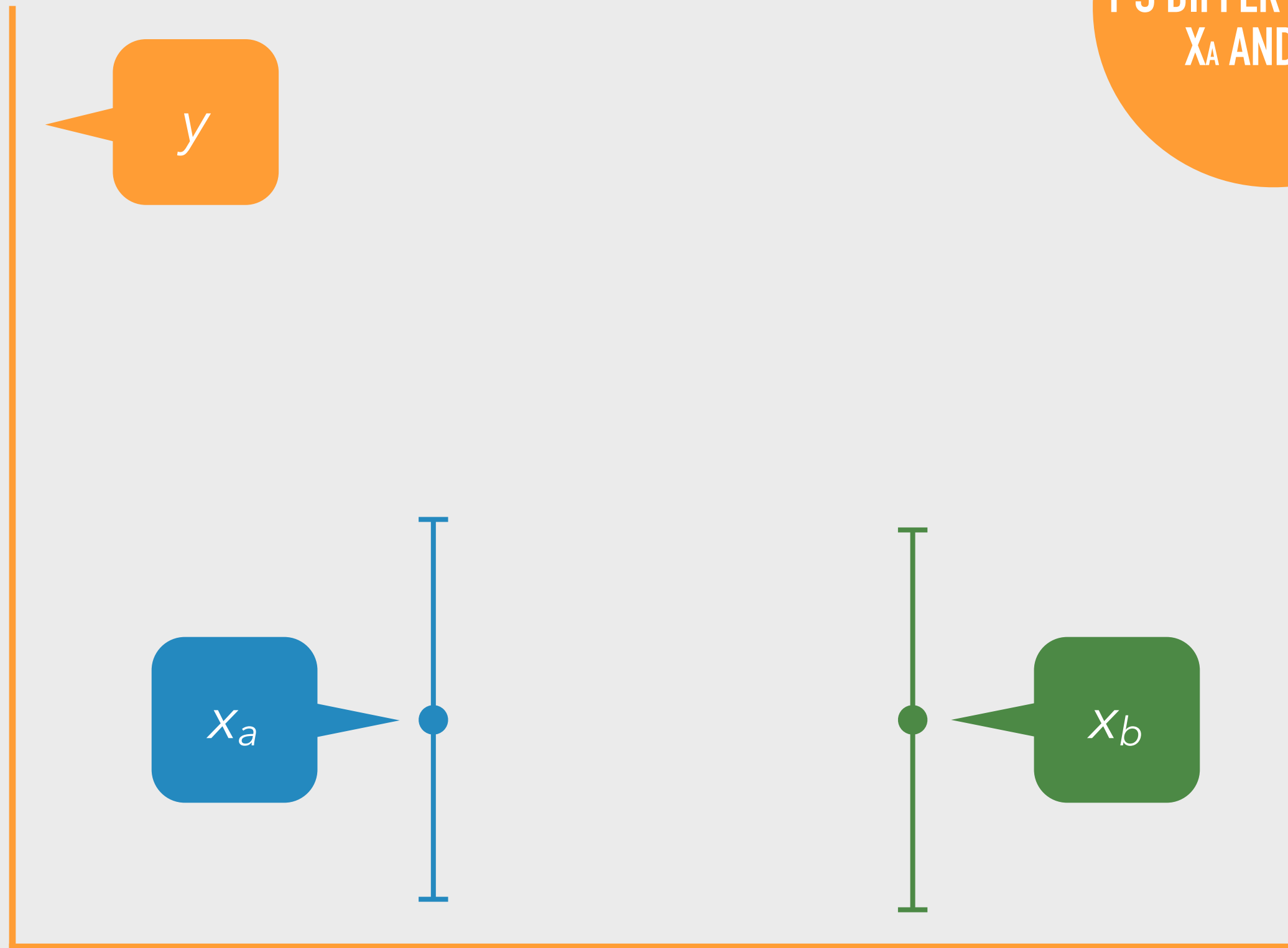
# INTERPRETATION

▸ The one-sample t-test (t=4.052,df=42, p<.001) suggests that these data are not representative of the population. The sample mean of 45 is significantly different from the population mean of 60.
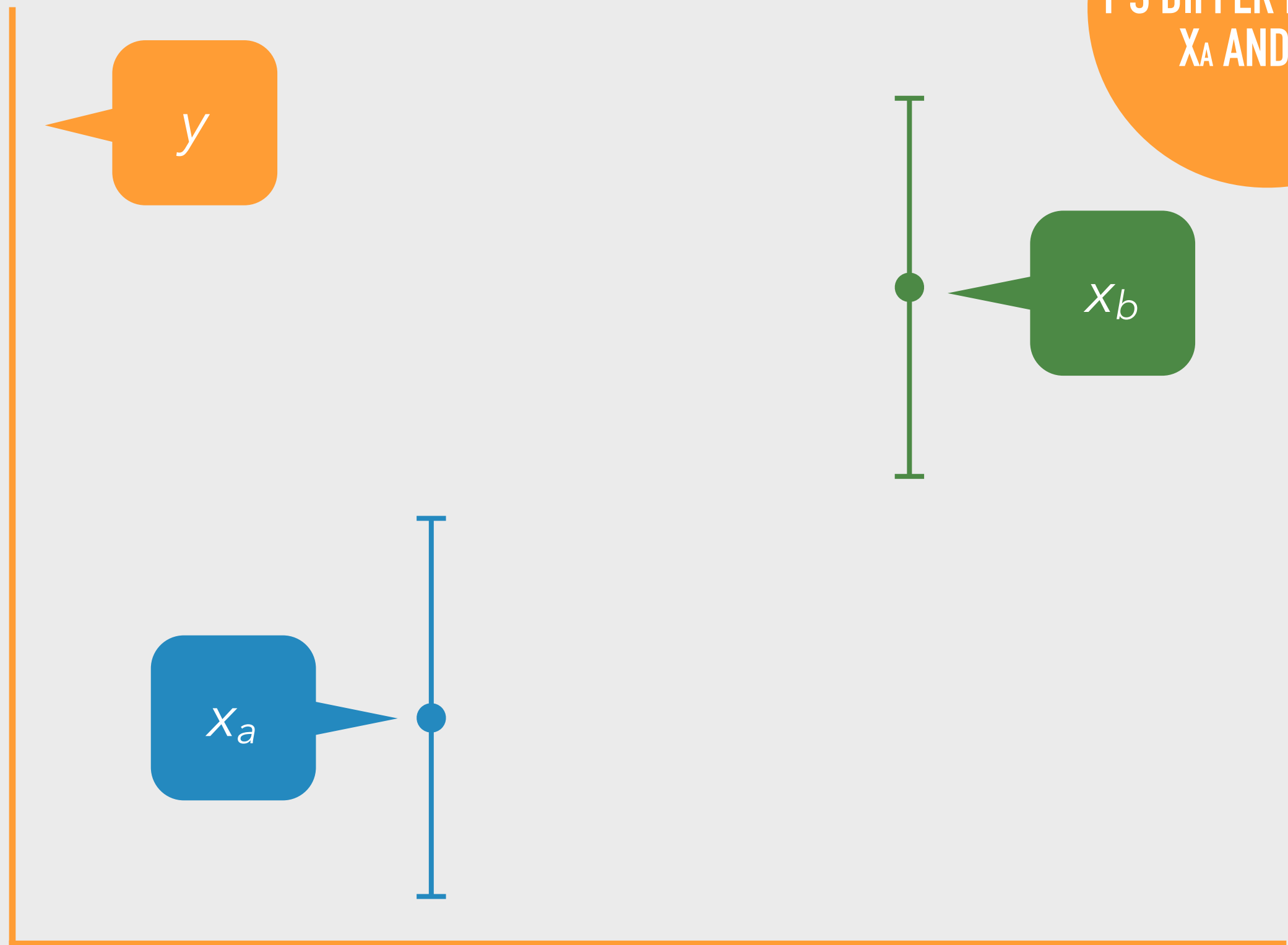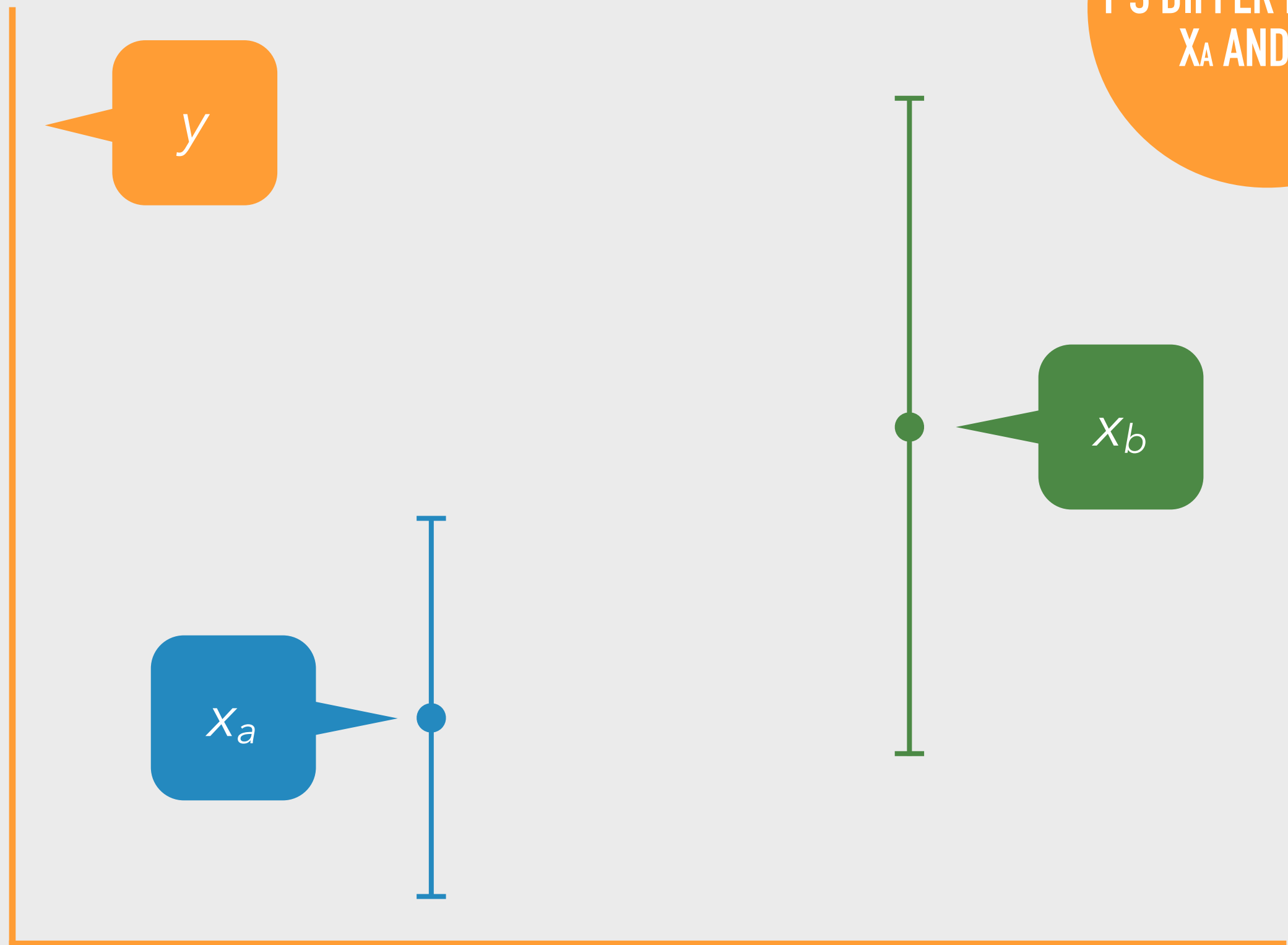
# 4 INDEPENDENT SAMPLES

# MODEL

# MODEL

# MODEL

# HYPOTHESES

▸ $H_0$ = there is no difference in the mean of *y* between $x_a$ and $x_b$

▸ $H_1$ = there is a difference in the mean of *y* between $x_a$ and $x_b$

# ASSUMPTIONS

▸ dependent variable ($y$) is continuous

▸ the distribution of $y$ is approximately normal

▸ independent variable is binary ($x_a$ and $x_b$)

▸ homogeneity of variance between $x_a$ and $x_b$

▸ observations are independent

▸ $v = n_a + n_b - 2$

# EQUATION ASSUMING HOMOGENEITY OF VARIANCE

$$t = \frac{\bar{X}_a - \bar{X}_b}{\sqrt{\frac{s_p^2}{n_a} + \frac{s_p^2}{n_b}}}$$

pooled variance

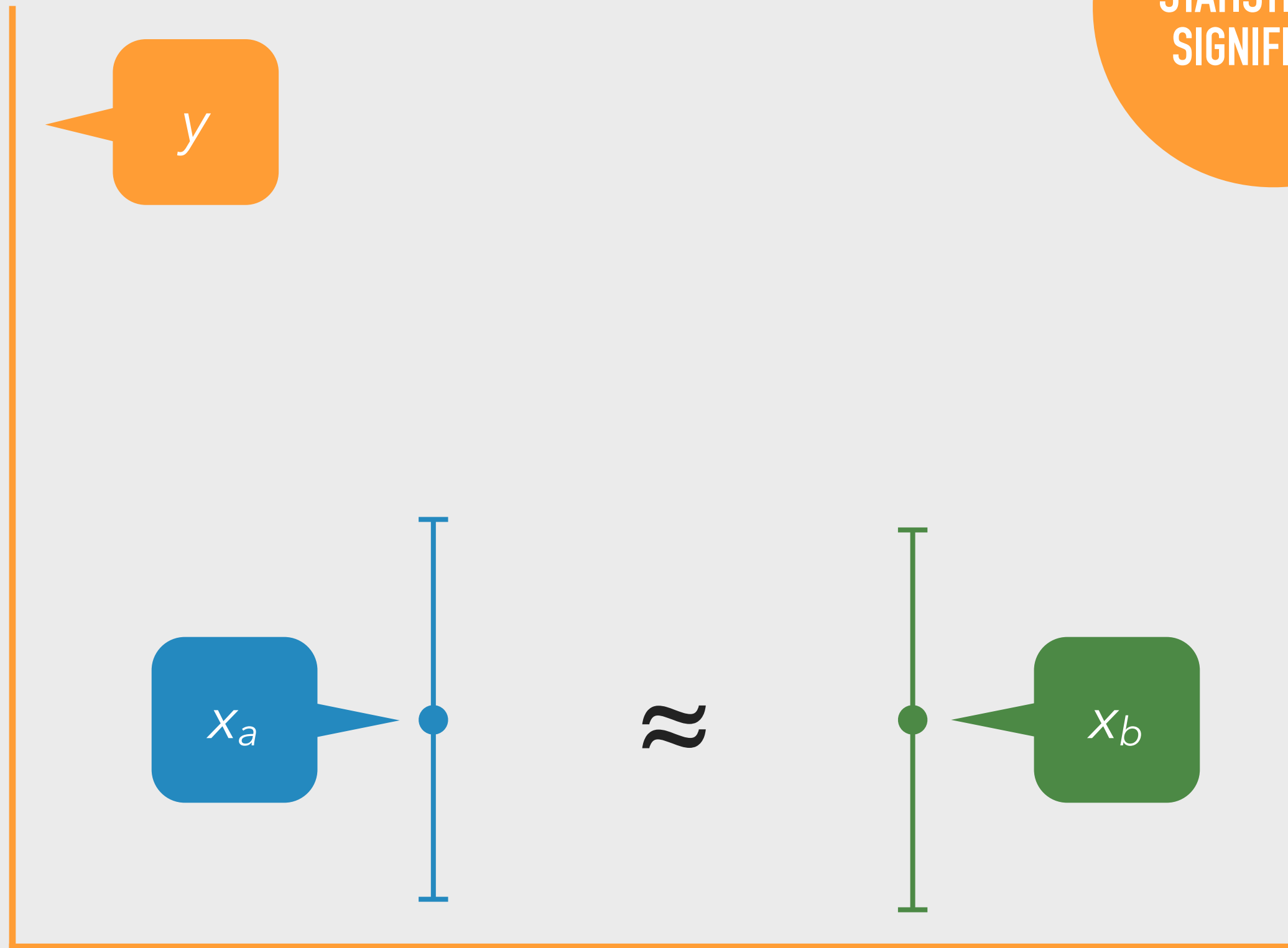# POOLED VARIANCE EQUATION

$$s_p^2 = \frac{(n_a - 1)\, s_a^2 + (n_b - 1)\, s_b^2}{n_a + n_b - 2}$$

# ASSUMPTIONS

▸ dependent variable ($y$) is continuous

▸ the distribution of $y$ is approximately normal

▸ independent variable is binary ($x_a$ and $x_b$)

▸ homogeneity of variance between $x_a$ and $x_b$

▸ observations are independent

▸ $v = n_a + n_b - 2$

# EQUATION IF HOMOGENEITY OF VARIANCE CANNOT BE ASSUMED

$$t = \frac{\bar{X}_a - \bar{X}_b}{\sqrt{\dfrac{s_a^2}{n_a} + \dfrac{s_b^2}{n_b}}}$$

variance values for each subgroup

# CAUTION! CAUTION! CAUTION!

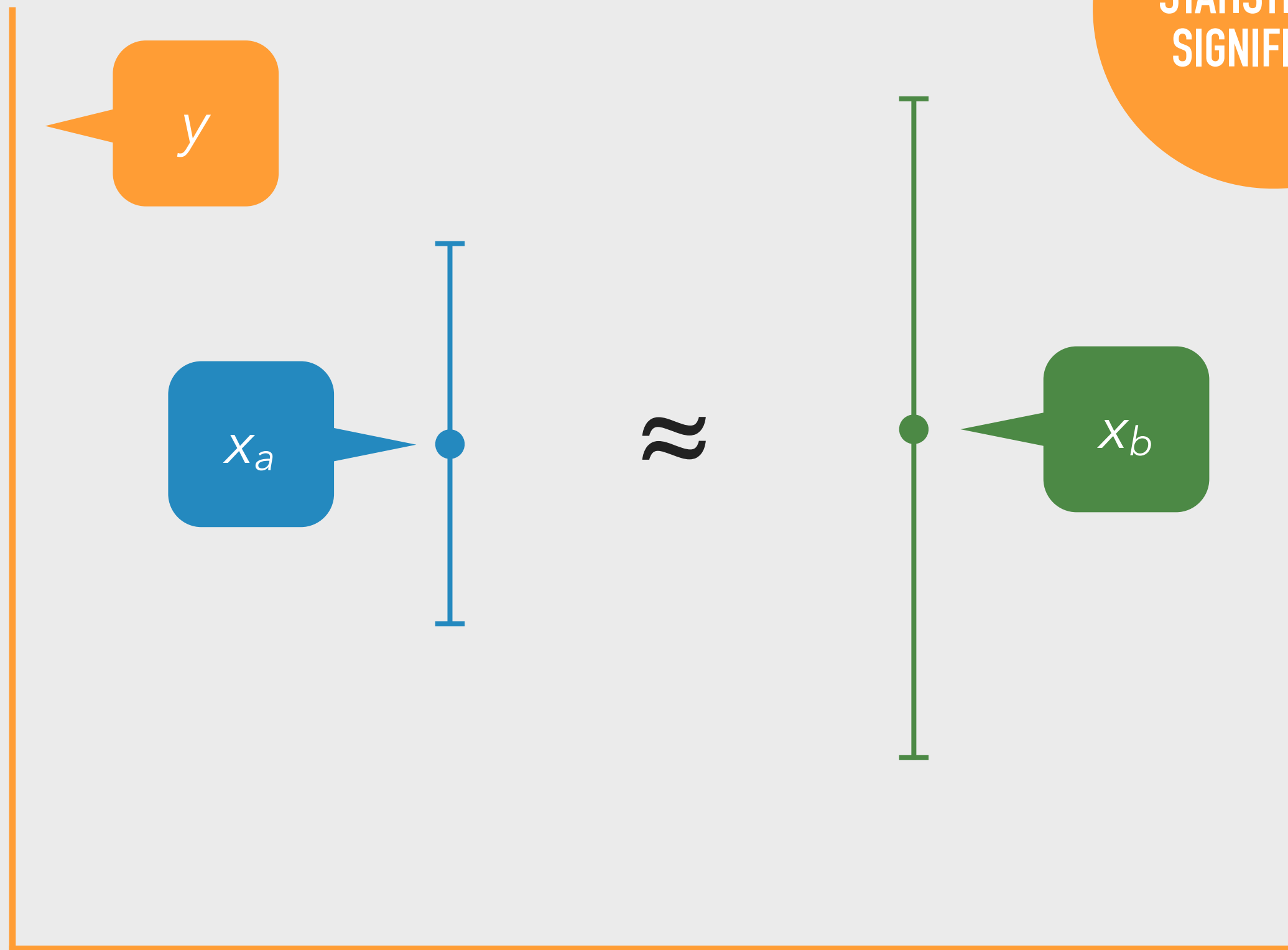$$t = \frac{\bar{X}_a - \bar{X}_b}{\sqrt{\frac{s_p^2}{n_a} + \frac{s_p^2}{n_b}}} \quad \neq \quad t = \frac{\bar{X}_a - \bar{X}_b}{\sqrt{\frac{s_a^2}{n_a} + \frac{s_b^2}{n_b}}}$$
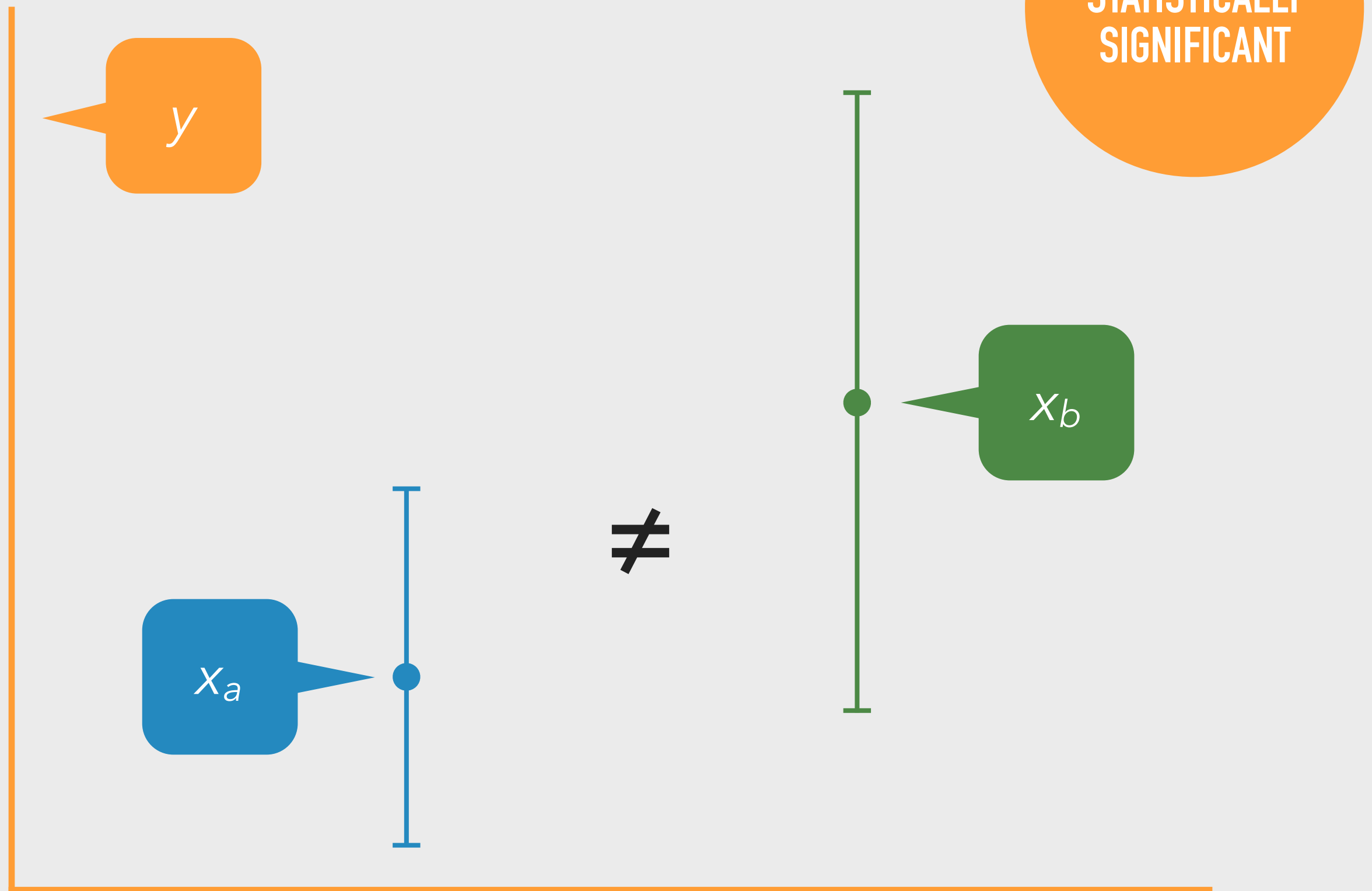
# WELCH'S CORRECTED DEGREES OF FREEDOM

$$v \approx \frac{\left( \frac{s_a^2}{n_a} + \frac{s_b^2}{n_b} \right)^2}{\frac{s_a^4}{(n_a^2)(n_a - 1)} + \frac{s_b^4}{(n_b^2)(n_b - 1)}}$$

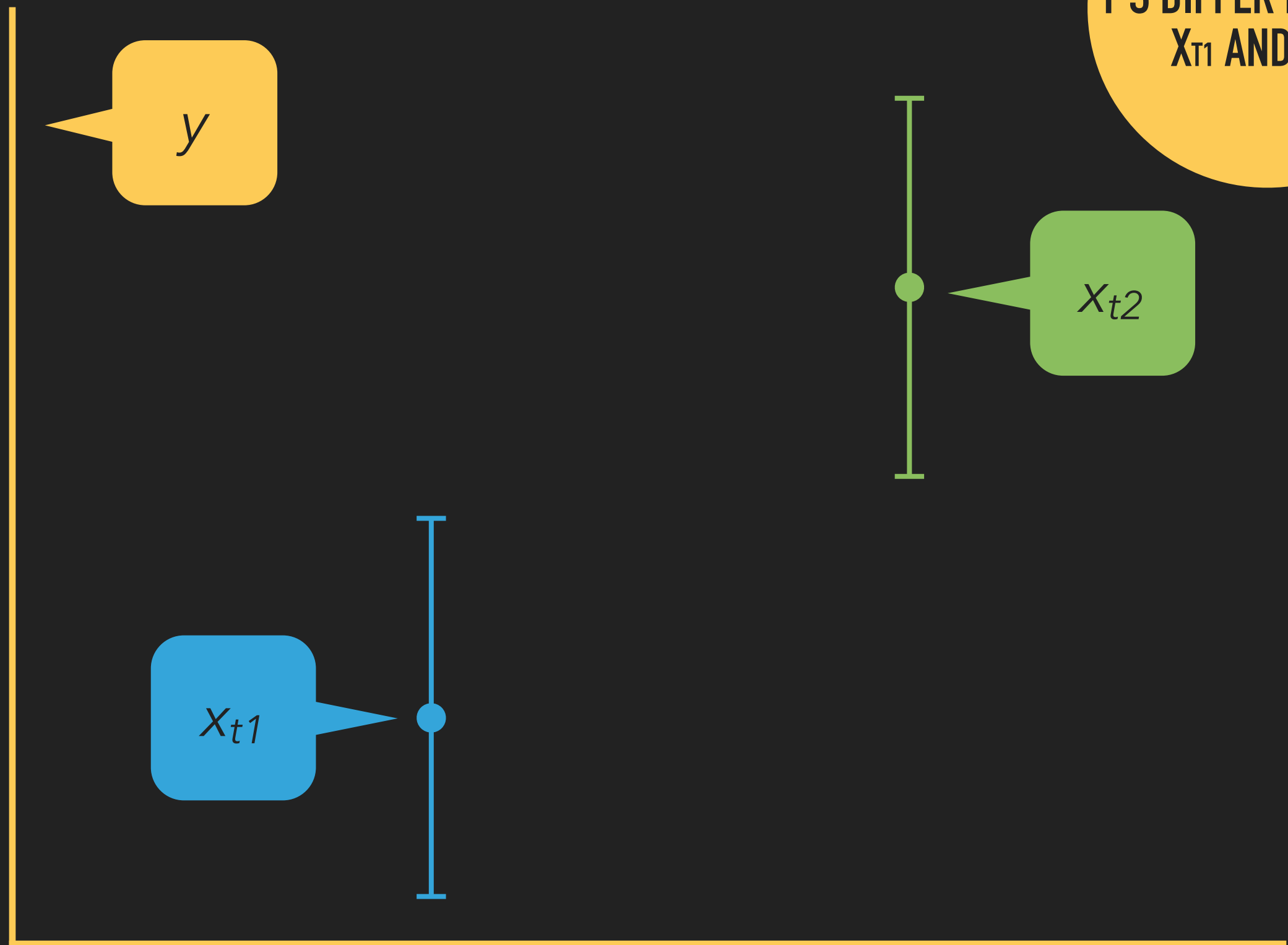# INTERPRETATION

▸ The independent t-test (t=4.052,df=42, p<.001) suggests that there is a significant difference in scores been men (mean of 20) and women (mean of 25). Results for women were found to be higher, on average, than results for men.

# 5 DEPENDENT SAMPLES

# MODEL

# MODEL

before

after

score

dependent variable

independent variable

# MODEL

# MODEL

# MODEL

# HYPOTHESES

▸ $H_0$ = there is no difference in the mean of *y* between $x_{t1}$ and $x_{t2}$

▸ $H_1$ = there is a difference in the mean of *y* between $x_{t1}$ and $x_{t2}$

# HYPOTHESES

▸ $H_0$ = there is no difference in the mean of *y* between $x_{g1}$ and $x_{g2}$

▸ $H_1$ = there is a difference in the mean of *y* between $x_{g1}$ and $x_{g2}$

# ASSUMPTIONS

▸ dependent variable ($y$) is continuous

▸ independent variable is binary ($x_{g1}$ and $x_{g2}$)
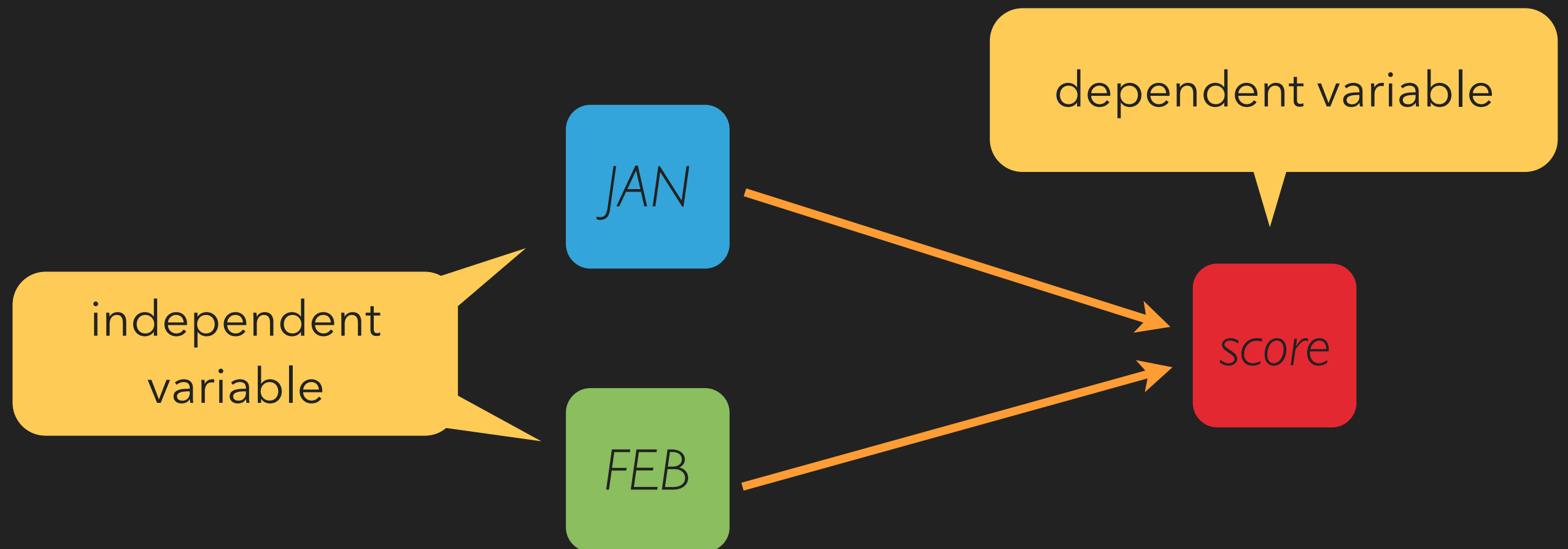
▸ the distribution of the differences between $x_{g1}$ and $x_{g2}$ is normally distributed

▸ scores are dependent

# EQUATION

$$t = \frac{\bar{d}}{\sqrt{\frac{s_d^2}{n}}}$$

mean of difference between groups

variance of difference between groups

# INTERPRETATION
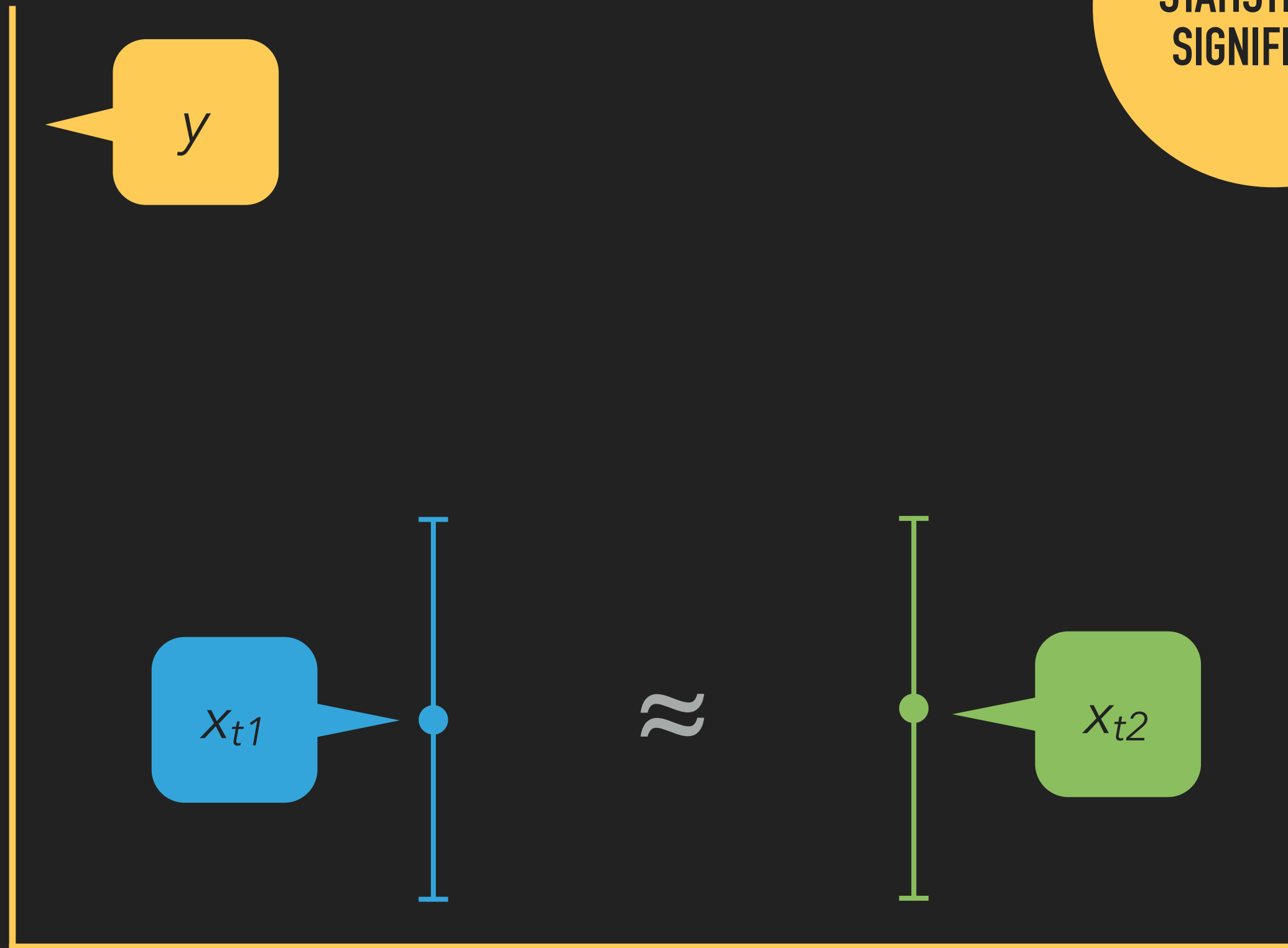
▸ The dependent t-test (t=4.052,df=42, p<.001) suggests that there is a significant difference in scores been the pre-test (mean of 20) and the post-test (mean of 25). Post-test results were found to be higher, on average, than pre-test results.

# LONG DATA

| participant | score | timePoint |
|---|---|---|
| jane | 10 | before |
| jane | 12 | after |
| john | 15 | before |
| john | 14 | after |

# WIDE DATA

| participant | score1 | score2 |
|---|---|---|
| jane | 10 | 12 |
| john | 15 | 14 |
| joe | 12 | 12 |
| jessica | 8 | 11 |

# RESHAPING DATA

| participant | score | timePoint |
|---|---|---|
| jane | 10 | before |
| jane | 12 | after |
| john | 15 | before |
| john | 14 | after |
| joe | 12 | before |
| joe | 12 | after |
| jessica | 8 | before |
| jessica | 11 | after |

↔

| participant | score1 | score2 |
|---|---|---|
| jane | 10 | 12 |
| john | 15 | 14 |
| joe | 12 | 12 |
| jessica | 8 | 11 |

# 6 EFFECT SIZES

Statistical Significance ≠ Real World Significance
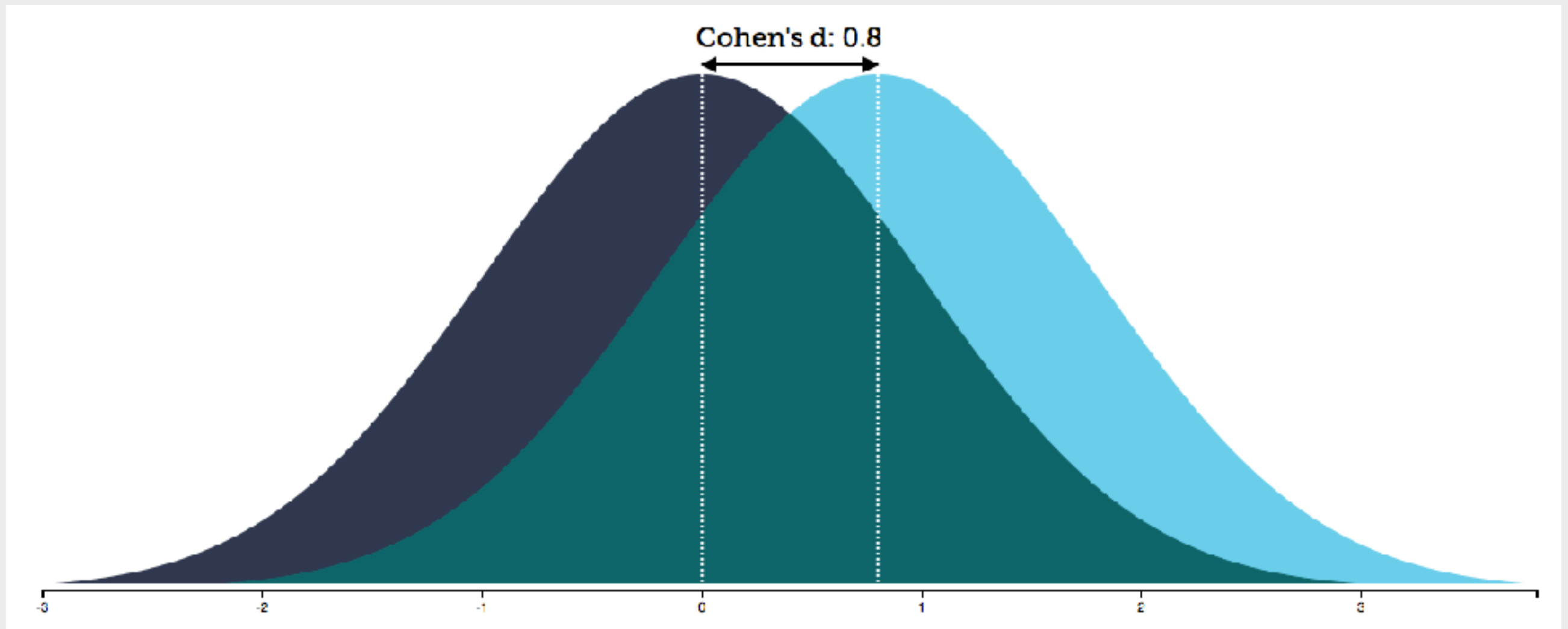
# COHEN'S D INTERPRETATION

# COHEN'S D INTERPRETATION

# COHEN'S D INTERPRETATION

# COHEN'S D EQUATION

$$d = \frac{\bar{x}_a - \bar{x}_b}{\sqrt{\frac{(n_a - 1)s_a^2 + (n_b - 1)s_b^2}{n_a + n_b - 2}}}$$

pooled variance

# COHEN'S D EQUATION

$$d = \frac{M_t - M_c}{\sqrt{\frac{(n_t - 1)s_t^2 + (n_c - 1)s_c^2}{n_t + n_c - 2}}}$$

# COHEN'S D EQUATION SIMPLIFIED

$n_a = n_b$

$$d = \frac{2t}{\sqrt{v}}$$

$n_a \neq n_b$

$$d = \frac{t\,(n_t + n_c)}{\sqrt{v}\,(\sqrt{n_t} + n_c)}$$

# DOCUMENT DETAILS

Document produced by Christopher Prener, Ph.D for the Saint Louis University course SOC 5050: QUANTITATIVE ANALYSIS - APPLIED INFERENTIAL STATISTICS. See the course wiki and the repository README.md file for additional details.