

CSCI 3022

intro to data science with probability & statistics

Lecture 19
March 19, 2018

Small sample size hypothesis testing

CSCI 3022

intro to data science with probability & statistics

Lecture 19
March 19, 2018

Small sample size hypothesis testing

(Sam Wang)

Stuff & Things

- HW5 posted tonight. Due the Friday *after* Spring Break.
- Dan's OH cancelled this Weds & Fri.

Previously on CSCI 3022

- Statistical inference for population mean **when data is normal** and n is large and...

- σ is known:

$$\left(\frac{\bar{x} - \mu}{\sigma / \sqrt{n}} \right) \sim N(0,1)$$

↑

- σ is unknown:

$$\left(\frac{\bar{x} - \mu}{s / \sqrt{n}} \right) \sim N(0,1)$$

↑
"empirical std. dev."

"z tests"

$$\underbrace{\bar{x}}_{\text{center}} \pm \underbrace{z_{\alpha/n} \sigma / \sqrt{n}}_{\text{window}}$$

Previously on CSCI 3022

- Statistical inference for population mean **when data is NOT normal** and n is large and...

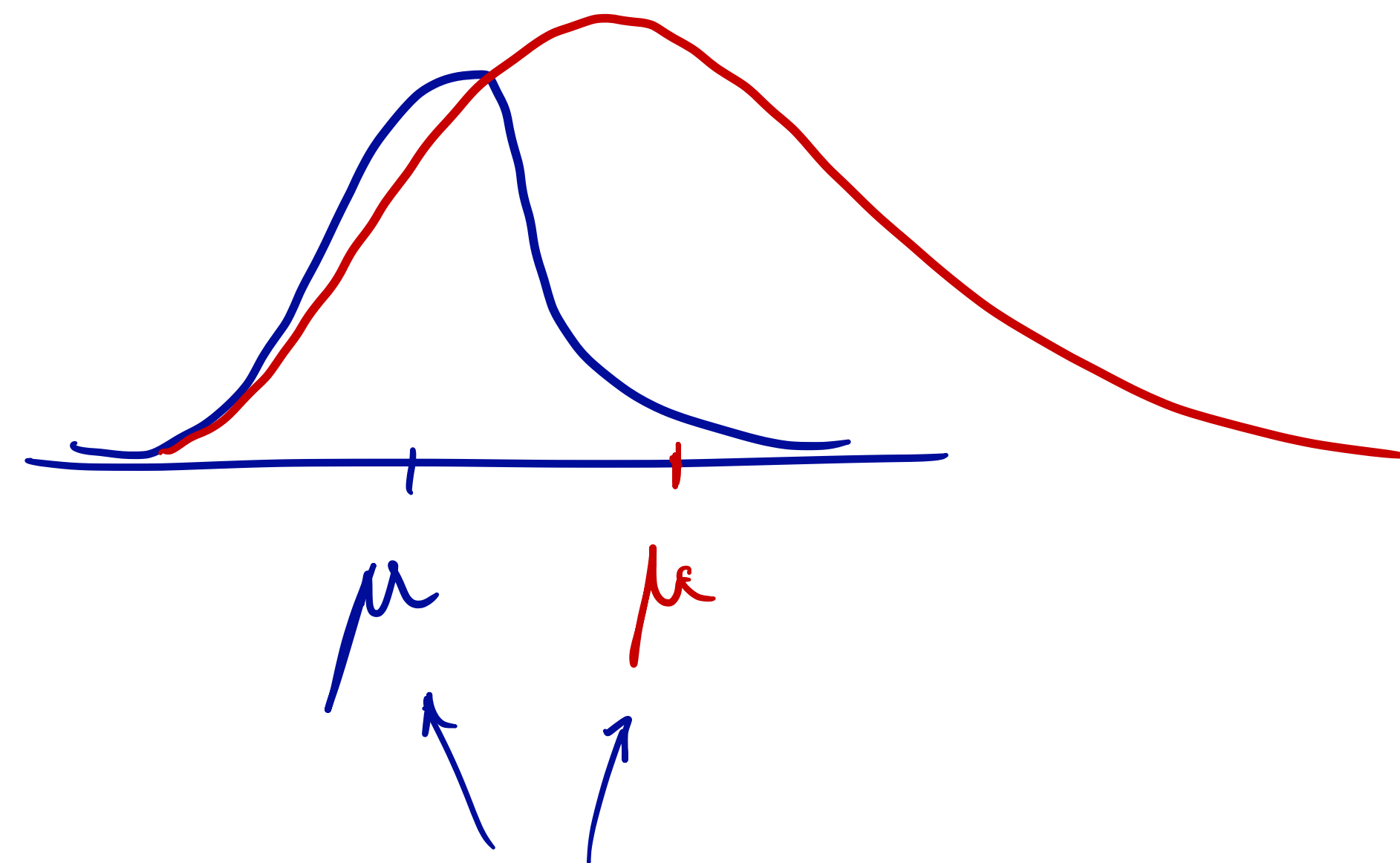
- σ is known:

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

- σ is unknown:

$$\left(\frac{\bar{X} - \mu}{s / \sqrt{n}} \right) \sim N(0, 1)$$

"Thanks, CLT!"



Previously on CSCI 3022

- Statistical inference for population mean **when data is normal** and n is small and...

$n < 30$



- σ is known:









$$\underbrace{\left(\frac{\bar{X} - \mu}{\sigma / \sqrt{n}} \right)}_{\text{???}} \sim N(0,1)$$

- σ is unknown:

???

The story so far, for means

- Thus far, we've talked about Hypothesis Testing & Confidence Intervals for the mean of a population in the following cases:

	"n is large" $n \geq 30$	"n is small" $n < 30$
Normal Data / Known σ		
Normal Data / Unknown σ ^{use s}		
Non-Normal Data / Known σ		
Non-Normal Data / Unknown σ		


 - z-test

 - t-test (TODAY!)

 Bootstrap
(after Spring Break)

Small-sample tests

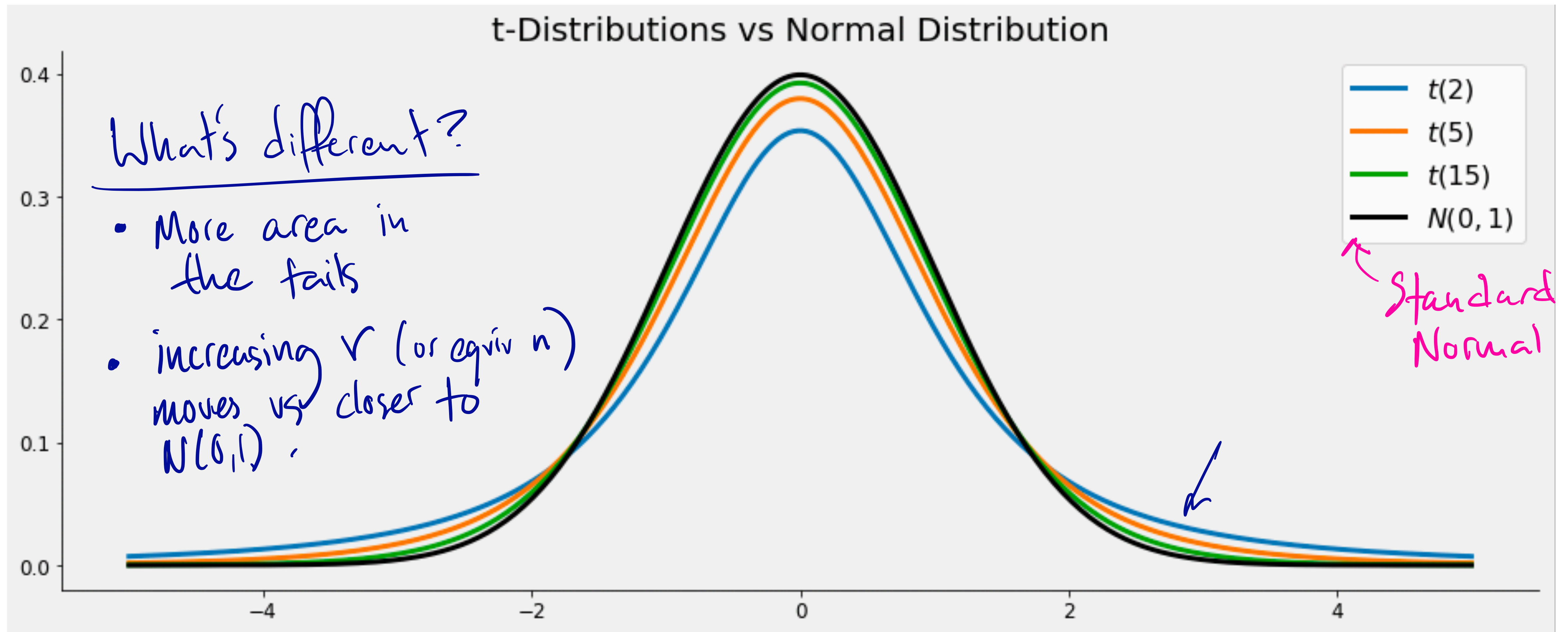
- When n is small we cannot invoke the Central Limit Theorem ☹
- When n is small and the variance is unknown we need to do something else ...
- When \bar{X} is the sample mean of a random sample of size n from a normal distribution with mean μ , the random variable

$$\left(\frac{\bar{X} - \mu}{s / \sqrt{n}} \right)$$


follows a probability distribution called a **t-Distribution** with parameter $\nu = n - 1$ degrees of freedom.

The t-Distribution

- The following figure shows the pdf of some members of the family of t-Distributions



- What do you notice about these t-Distributions, compared with the Standard Normal curve?

Properties of t-Distributions

- Let t_ν denote the t-Distribution with parameter ν degrees of freedom
- Each t_ν -curve is bell-shaped and centered at 0
- Each t_ν -curve is more spread out than the standard normal distribution
- As ν increases, the spread of the corresponding t_ν -curve decreases
- As $\nu \rightarrow \infty$ the sequence of t_ν -curves approaches the standard normal curve

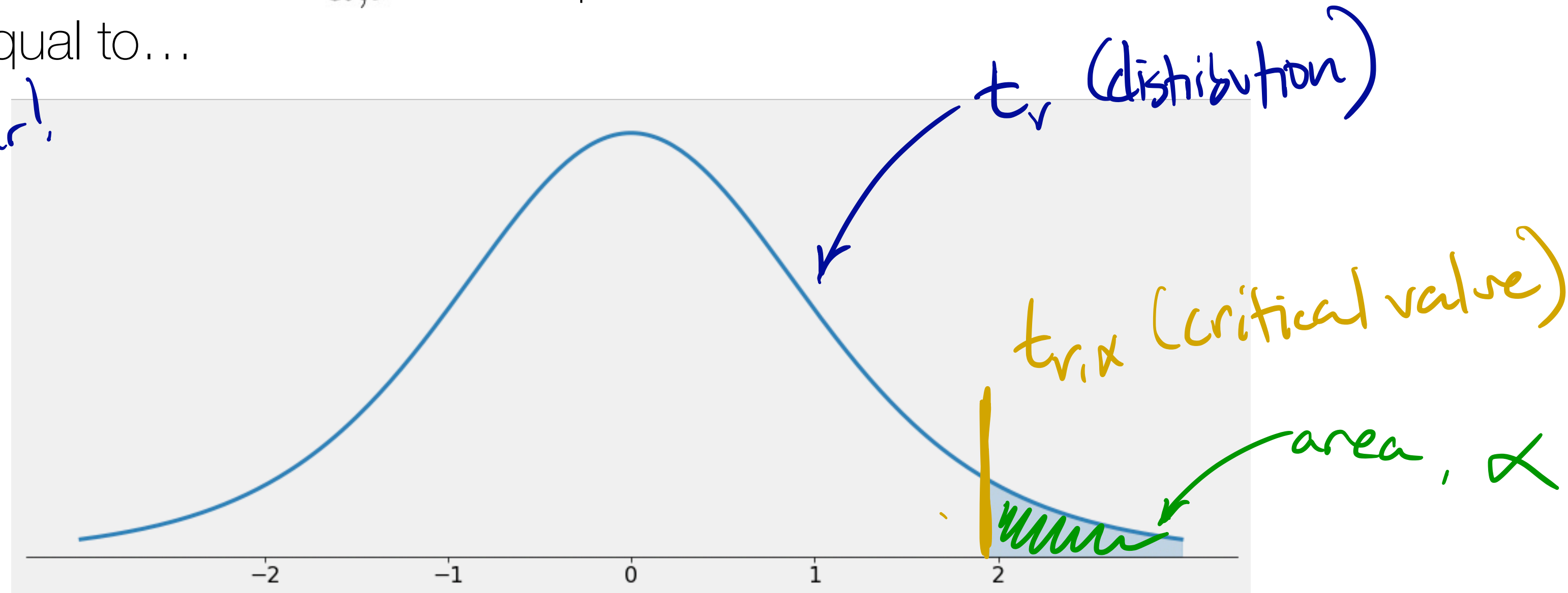
$\nu = n - 1$

Aside:
 $\backslash nu$ in
LaTeX.

The t-critical value

- We can extend all of our inferential mechanics to the small-sample case by introducing the so-called t-critical value, which we denote $t_{\alpha, \nu}$
- **Definition:** the t-critical value $t_{\alpha, \nu}$ is the point such that the area under the t_{ν} -curve to the right of $t_{\alpha, \nu}$ is equal to...

This should look familiar!
Very similar to our old friend, the z-test,
and using Z-critical values



- Example: $t_{0.05, 6}$ is the t-critical value that captures the upper-tail area of 0.05 under the t curve with 6 degrees of freedom.

The t-confidence interval for the mean

- Let \bar{x} and s be the sample mean and sample standard deviation computed from the results of a random sample with of size n from a normal population with mean μ .

- Then a $100(1 - \alpha)\%$ t-confidence interval for the mean μ is given by:

$$\left[\bar{x} - t_{\alpha/2, n-1} \frac{s}{\sqrt{n}}, \bar{x} + t_{\alpha/2, n-1} \frac{s}{\sqrt{n}} \right]$$

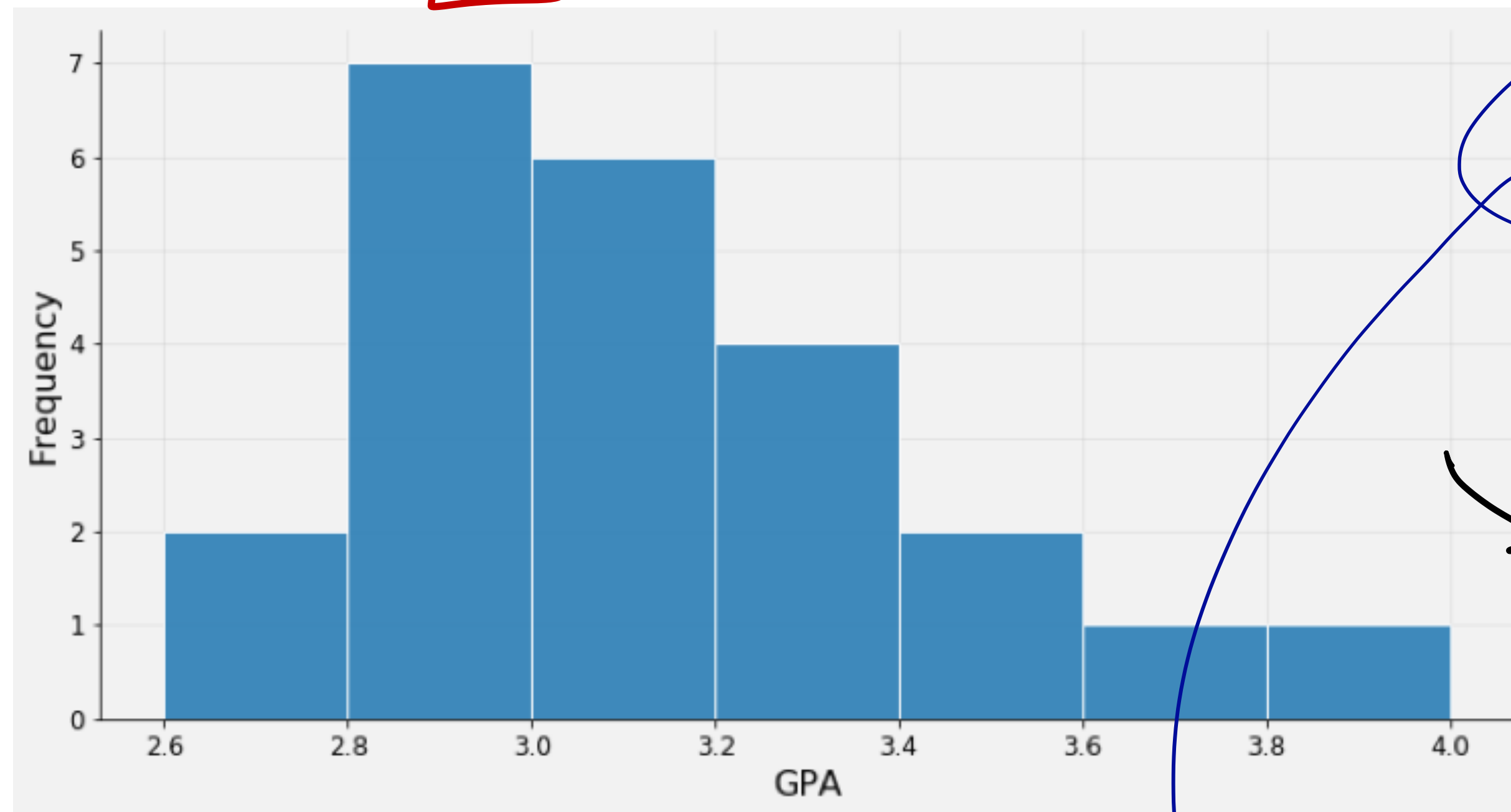
- Or more compactly:

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

CI

t-confidence interval example

- **Example:** Suppose the GPAs for 23 students have a histogram that looks as follows:



$$\begin{aligned}n &= 23 \\ \bar{x} &= 3.146 \\ s &= 0.308 \\ \alpha &= 0.1 \\ \alpha/2 &= 0.05\end{aligned}$$

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

$$t_{\alpha/2, n-1}$$

$$\begin{aligned}\text{Stats.t.ppf}(0.95, 23-1) \\ = 1.717\end{aligned}$$

\uparrow
 $n-1$

- The sample mean of the data is 3.146 and the sample standard deviation is 0.308. Find a 90% confidence interval for the mean GPA.

$$\begin{aligned}\alpha &= 0.1 \\ "(1-\alpha) \cdot 100" \% CI\end{aligned}$$

$$3.146 \pm 1.717 \cdot \frac{0.308}{\sqrt{23}}$$

$$\Rightarrow [3.033, 3.259]$$

The t-Test, Critical Regions and P-Values

$$H_0 : \theta = \theta_0$$

Alternative Hypothesis

$$H_1 : \theta > \theta_0$$

$$H_1 : \theta < \theta_0$$

$$H_1 : \theta \neq \theta_0$$

Critical Region Level α Test

$$t \geq t_{\alpha, \nu}$$

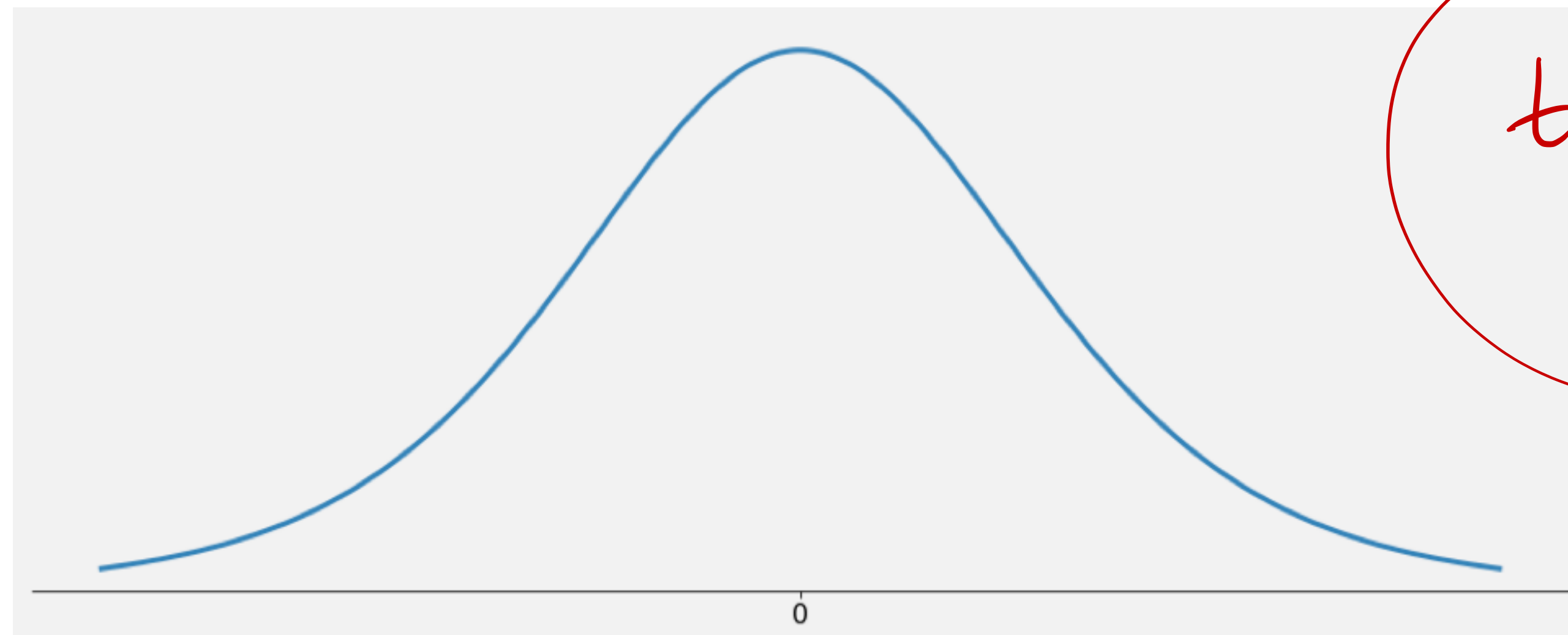
$$t \leq t_{\alpha, \nu}$$

$$(t \leq -t_{\alpha/2, \nu}) \text{ or } (t \geq t_{\alpha/2, \nu})$$

t test statistic
looks just like
z test statistics!

The only
difference
... is n is small
($n < 30$)

confidence
degrees of freedom



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

"standardized
statistic"

The t-Test, Critical Regions and P-Values

Alternative Hypothesis

P-Value Level α Test

$$H_1 : \theta > \theta_0$$

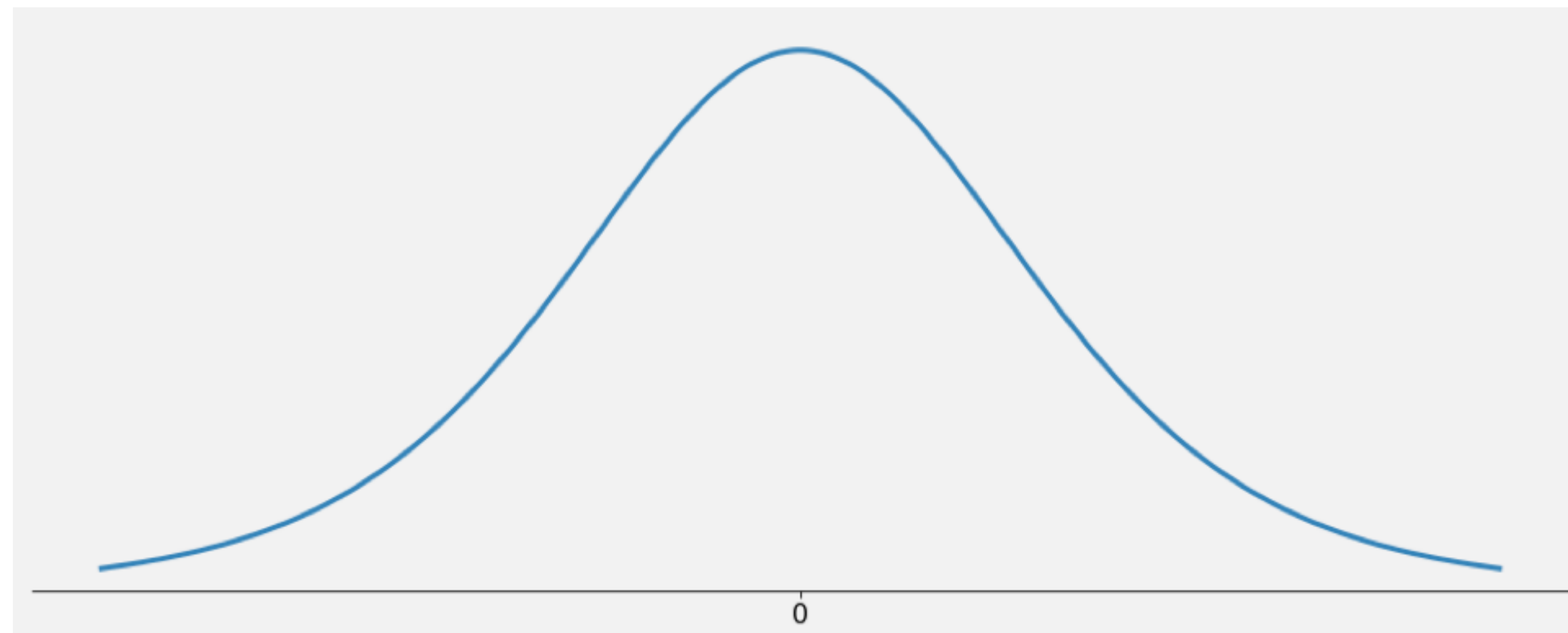
$$P(T \geq t \mid H_0) \leq \alpha$$

$$H_1 : \theta < \theta_0$$

$$P(T \leq t \mid H_0) \leq \alpha$$

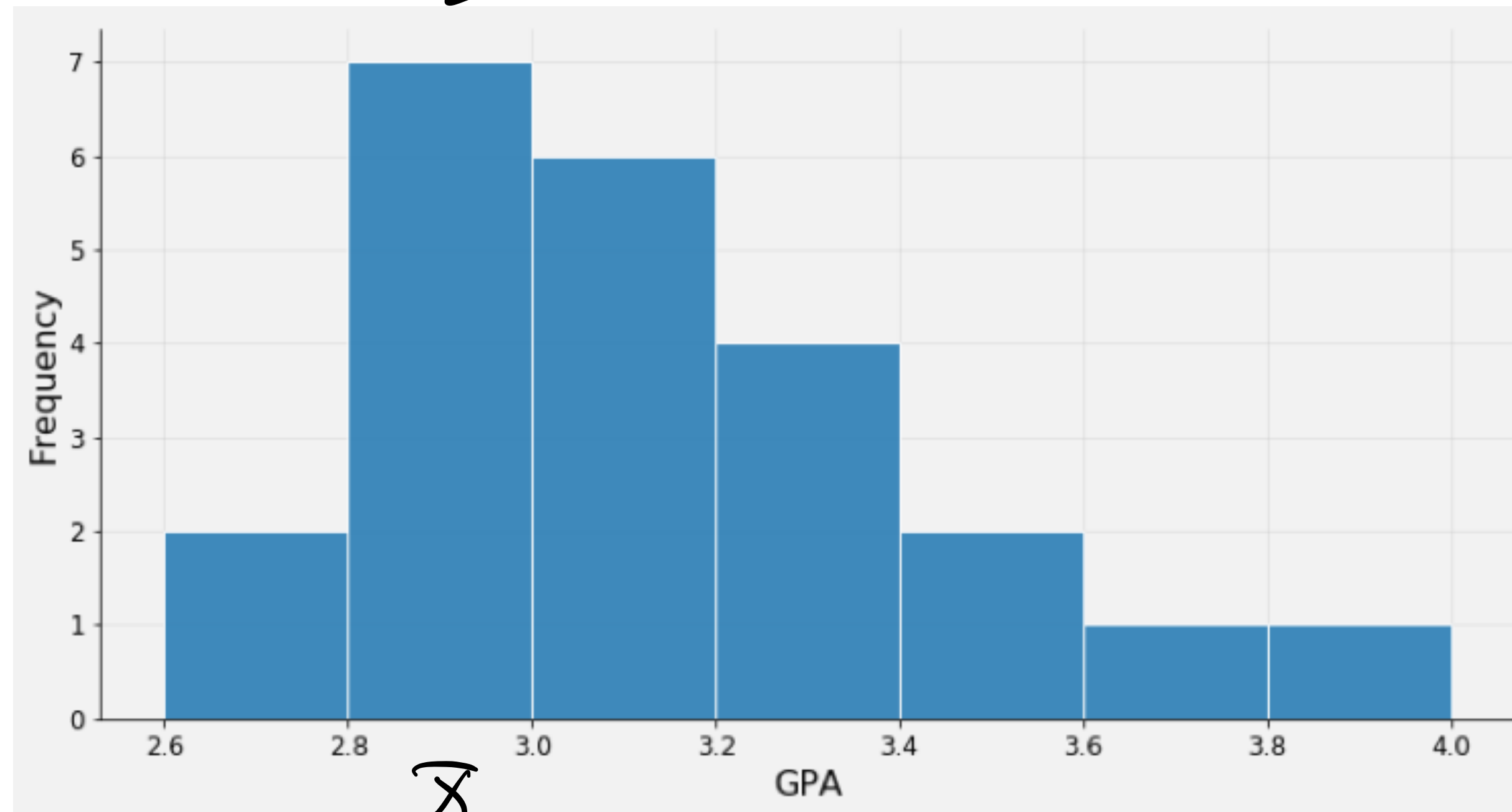
$$H_1 : \theta \neq \theta_0$$

$$2 \min \{P(T \leq t \mid H_0), P(T \geq t \mid H_0)\} \leq \alpha$$



t-Test example (p-value method)

- **Example:** Suppose the GPAs for 23 students have a histogram that looks as follows:



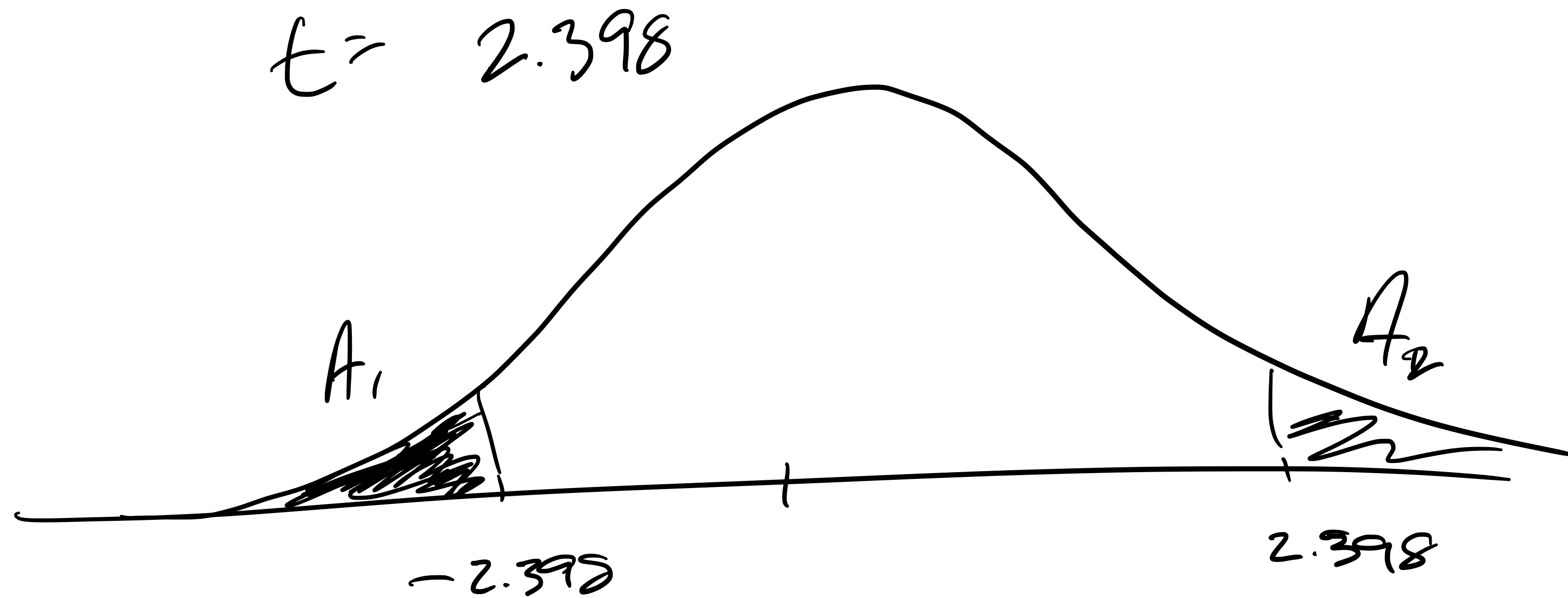
- The sample mean of the data is 3.146 and the sample standard deviation is 0.308. Determine if there is sufficient evidence to conclude at the 0.10 significance level that the mean GPA is not equal to 3.30.

$$H_1: \text{GPA} \neq 3.30$$

$$\alpha = 0.1$$

$$t = \frac{\bar{X} - \mu}{s / \sqrt{n}} = \frac{3.146 - 3.30}{0.308 / \sqrt{23}} = -2.398$$

t-Test example (p-value method)



$$2 \times \text{stats.t.cdf}(\underset{\substack{\uparrow \\ \text{test} \\ \text{Statistic}}}{-2.398}, \underset{\text{dof}}{22}) = \boxed{0.0254} < 0.10$$

$\underbrace{\hspace{10em}}_{\text{p-value}} \quad \alpha$