# intro to data science
# with probability & statistics

## CSCI 3022

Lecture 19
March 19, 2018

Small sample size hypothesis testing

Department of Computer Science
UNIVERSITY OF COLORADO **BOULDER**

Dan Larremore

# intro to data science
# with probability & statistics

# CSCI 3022

Lecture 19
March 19, 2018

Small sample size hypothesis testing

Department of Computer Science
UNIVERSITY OF COLORADO **BOULDER**

Dan Larremore

# 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…

- $\sigma$ is known:

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

"Z tests"

$$\bar{X} \pm Z_{\alpha/n} \, \sigma / \sqrt{n}$$

Center     Window

- $\sigma$ is unknown:

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

"empirical std. dev."

# Previously on CSCI 3022

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

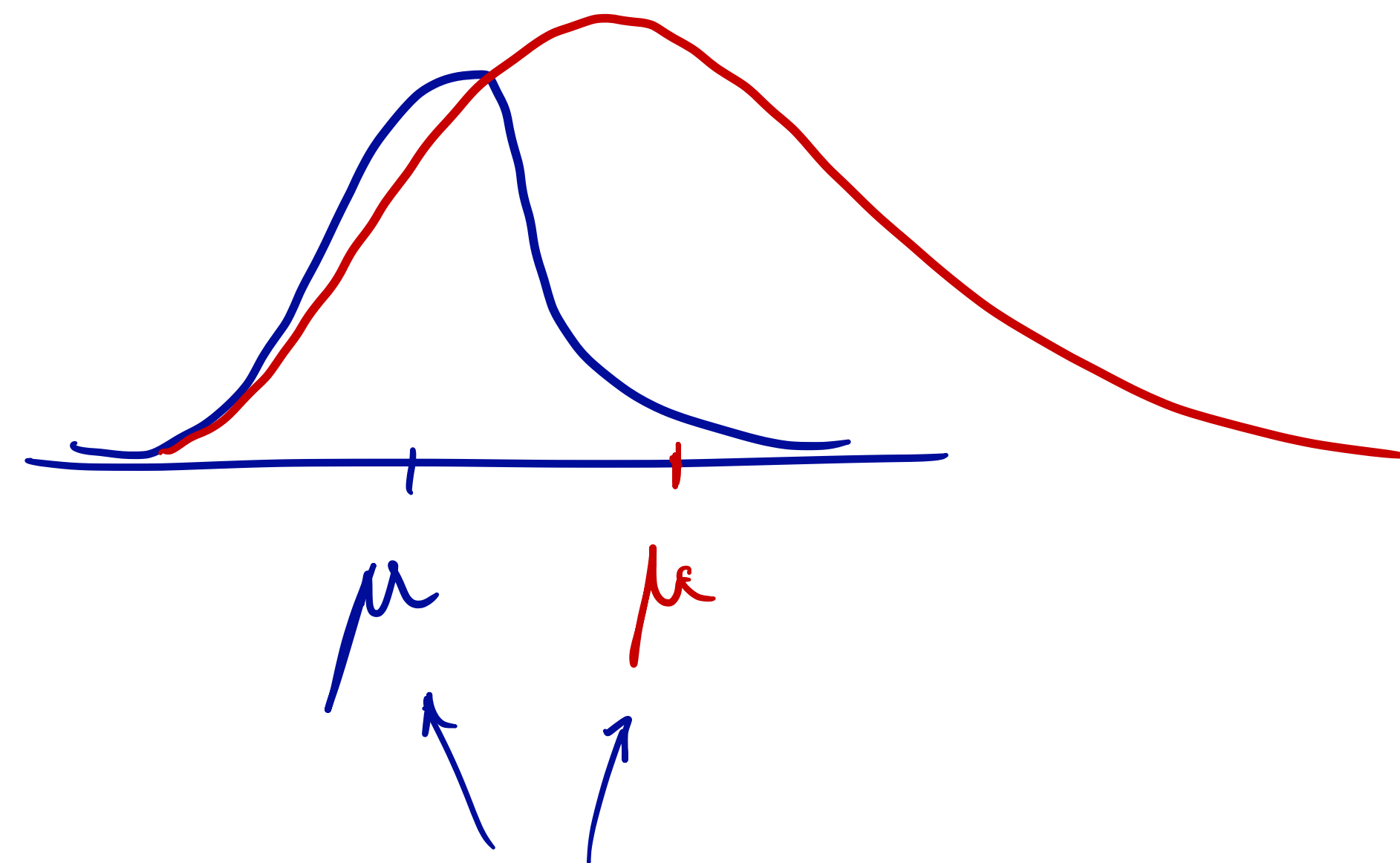- $\sigma$ is known:

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

"Thanks, CLT!"

- $\sigma$ is unknown:

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

$\mu$ $\mu$

# Previously on CSCI 3022

$n < 30$

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

- $\sigma$ is known:

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

- $\sigma$ 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 $\sigma$ | 〰️ (orange) | 〰️ (orange) |
| Normal Data / Unknown $\sigma$ — use $S$ | 〰️ (orange) | 〰️ (blue) |
| Non-Normal Data / Known $\sigma$ | 〰️ (orange) | 〰️ (green) |
| Non-Normal Data / Unknown $\sigma$ | 〰️ (orange) | 〰️ (green) |

〰️ (orange) — z-test
〰️ (blue) — t test (TODAY!)
〰️ (green) 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 $\mu$, 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

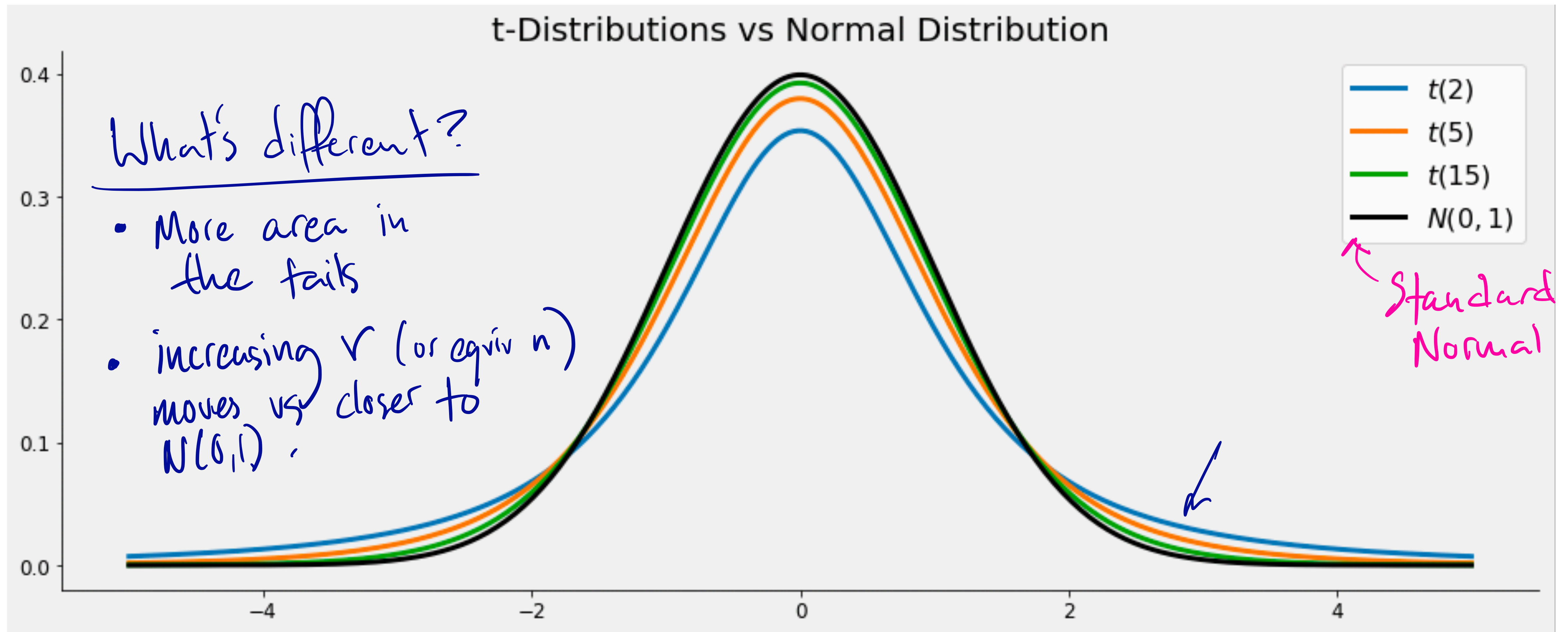- The following figure shows the pdf of some members of the family of t-Distributions



t-Distributions vs Normal Distribution

Legend:
- $t(2)$
- $t(5)$
- $t(15)$
- $N(0,1)$

Handwritten annotations:

What's different?

- More area in the tails
- increasing $\nu$ (or equiv $n$) moves us closer to $N(0,1)$.

↑ Standard Normal

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

# Properties of t-Distributions

- Let $t_\nu$ denote the t-Distribution with parameter $\nu$ degrees of freedom

- Each $t_\nu$ -curve is bell-shaped and centered at 0

- Each $t_\nu$-curve is more spread out than the standard normal distribution

- As $\nu$ increases, the spread of the corresponding $t_\nu$-curve decreases

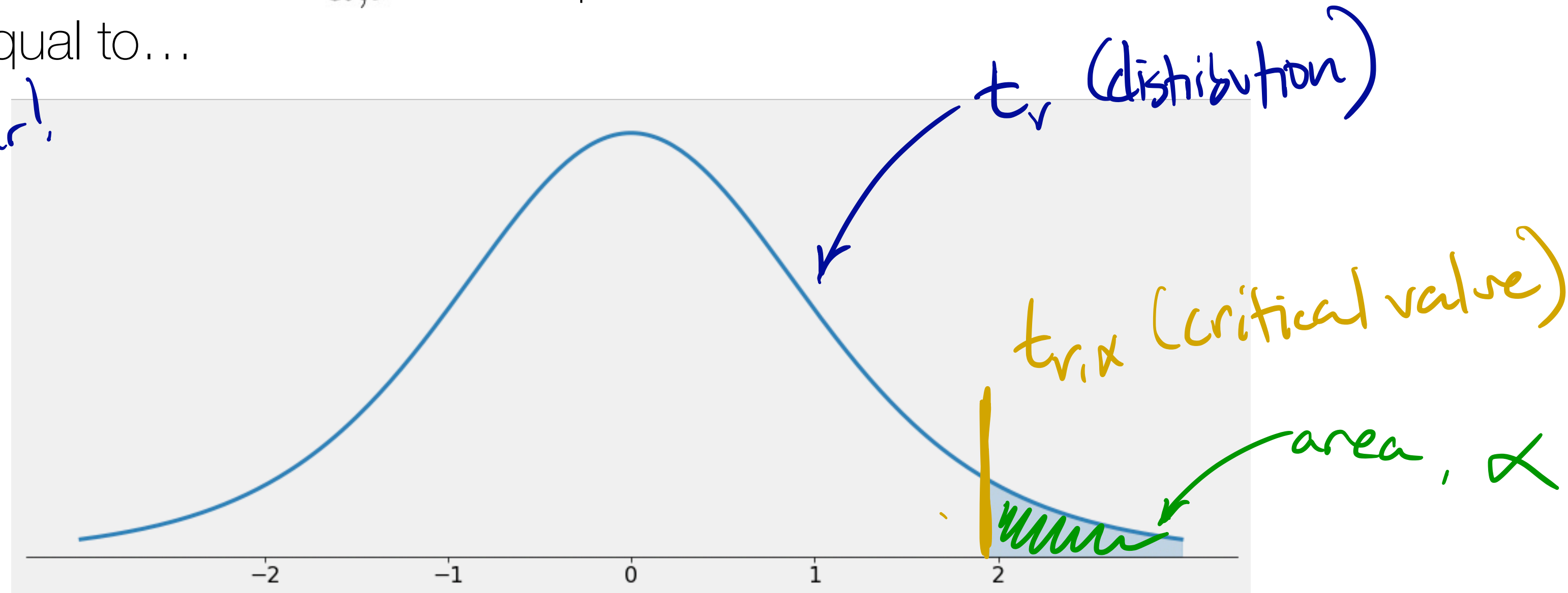- As $\nu \to \infty$ the sequence of $t_\nu$ -curves approaches the standard normal curve

$\nu = n-1$

$\left( \begin{array}{c} \text{Aside:} \\ \text{\textbackslash nu in} \\ \text{LaTeX.} \end{array} \right)$

# 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 $t_\nu$ -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 $z_{\alpha/2}$ critical values*

$t_\nu$ (distribution)

$t_{\nu,\alpha}$ (critical value)

area, $\alpha$

- 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 $\mu$.

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

$$\left[ \bar{X} - t_{\alpha/2, n-1} \, s/\sqrt{n}, \; \bar{X} + t_{\alpha/2, n-1} \, s/\sqrt{n} \right]$$
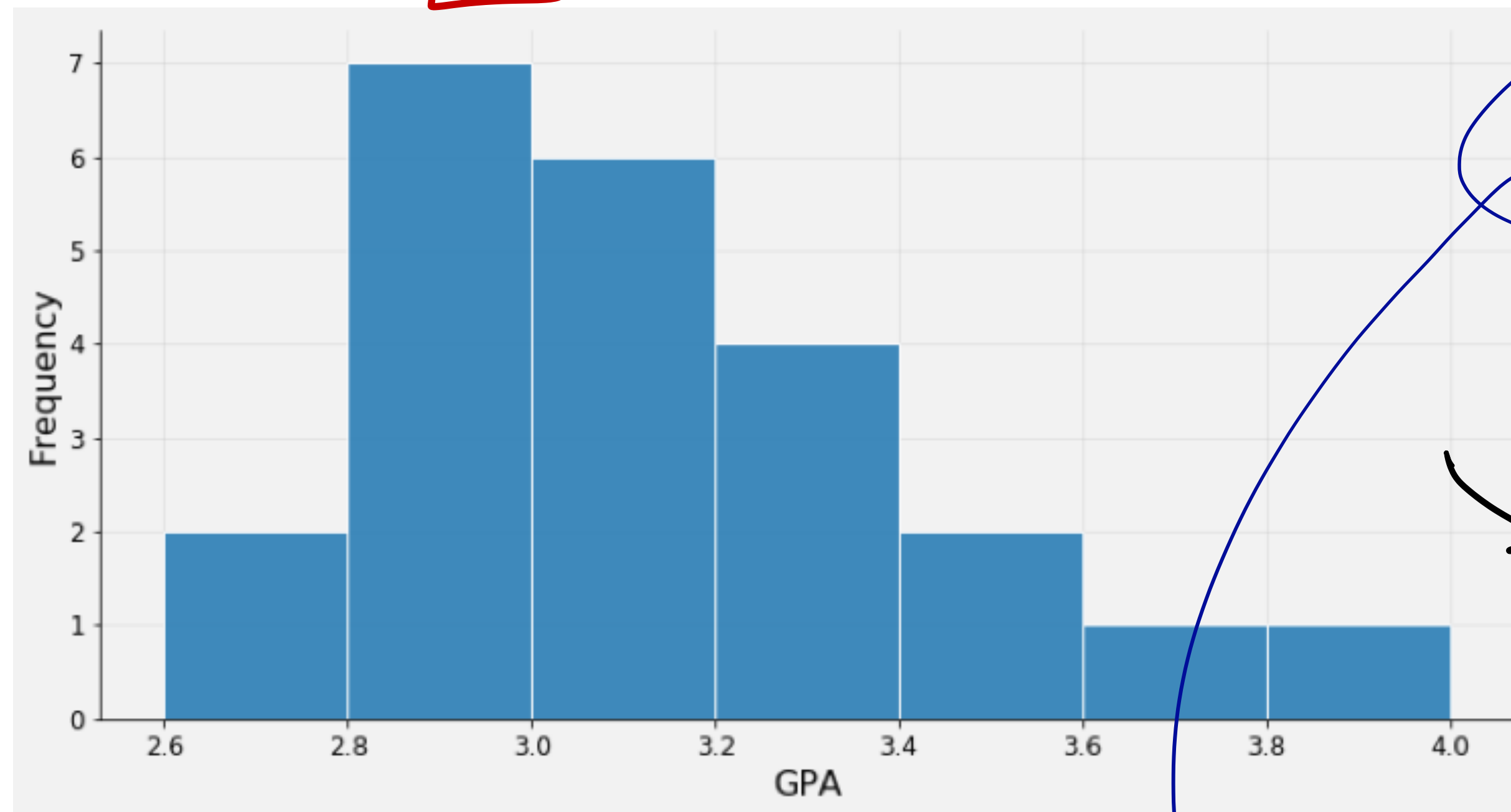
- Or more compactly:

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

CI

# t-confidence interval example

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

$n = 23$
$\bar{x} = 3.146$
$s = 0.308$
$\alpha = 0.1$

$\alpha/2 = 0.05$



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

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

$\text{stats.t.ppf}(0.95, 23-1)$
$= 1.717$
$n-1$

- 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.

$\rightarrow \alpha = 0.1$

"$(1-\alpha) \cdot 100$" % CI

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

$\Rightarrow [3.033, 3.259]$

# The t-Test, Critical Regions and P-Values

$$H_0 : \quad \theta = \theta_0$$

**Alternative Hypothesis**                          **Critical Region Level** $\alpha$ **Test**

t test statistic looks just like z test statistics!

$$H_1 : \quad \theta > \theta_0 \qquad\qquad t \geq t_{\alpha,\nu}$$

confidence

degrees of freedom

$$H_1 : \quad \theta < \theta_0 \qquad\qquad t \leq t_{\alpha,\nu}$$

$$H_1 : \quad \theta \neq \theta_0 \qquad (t \leq -t_{\alpha/2,\nu}) \text{ or } (t \geq t_{\alpha/2,\nu})$$

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

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

"standardized statistic"

# The t-Test, Critical Regions and P-Values

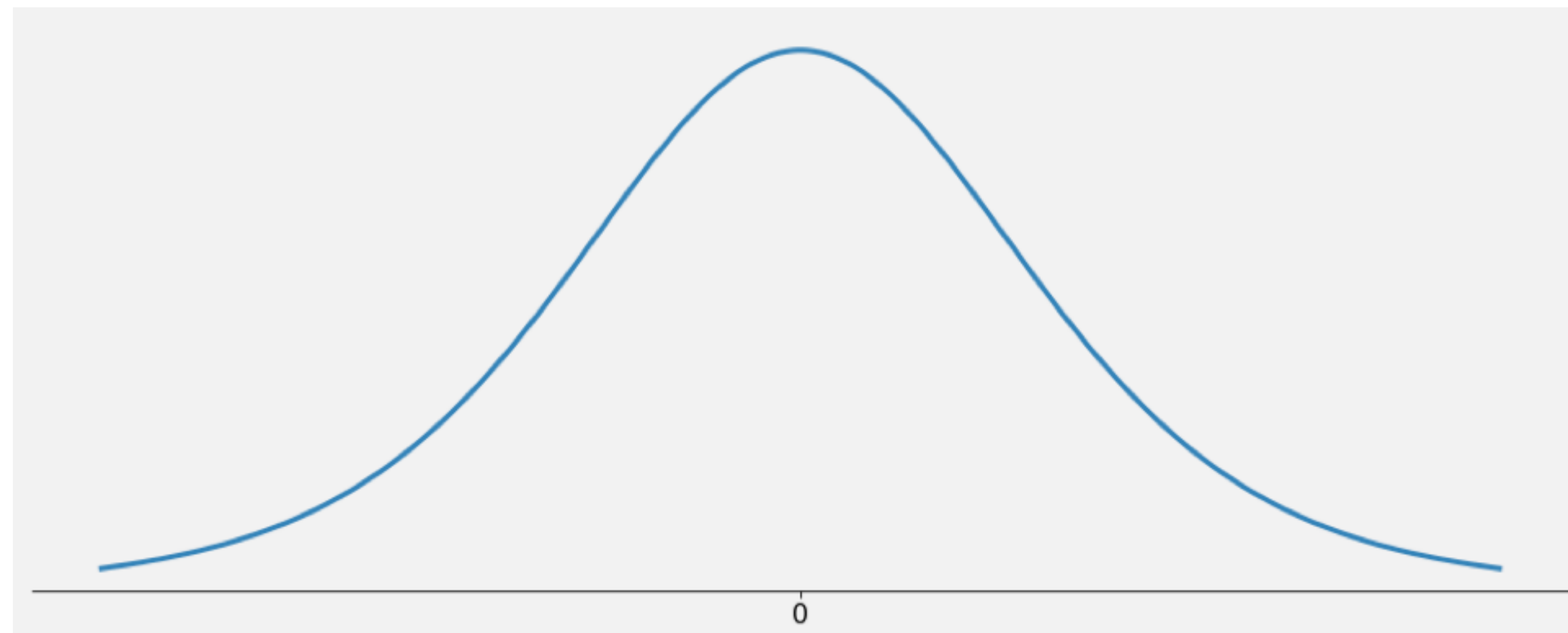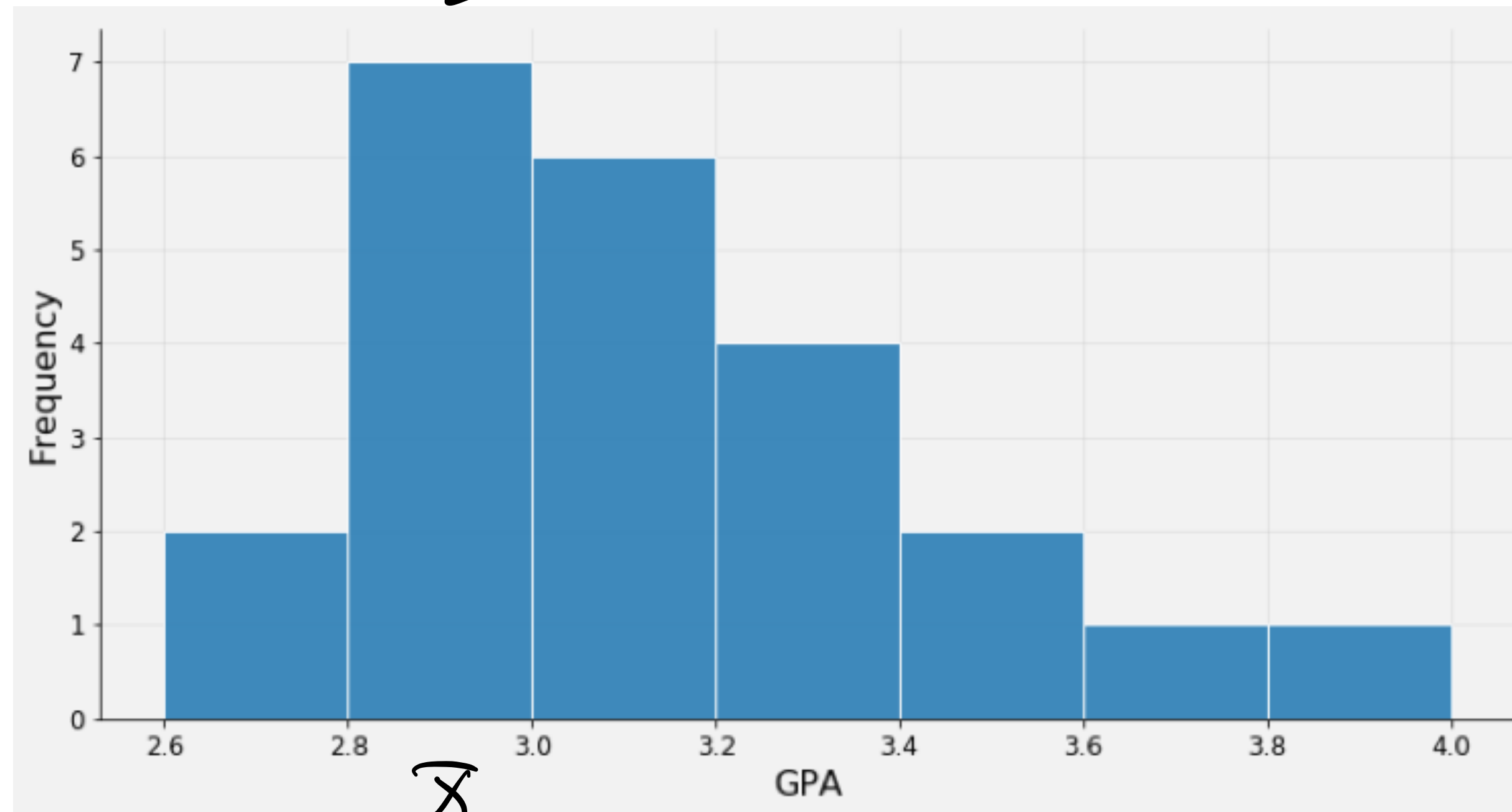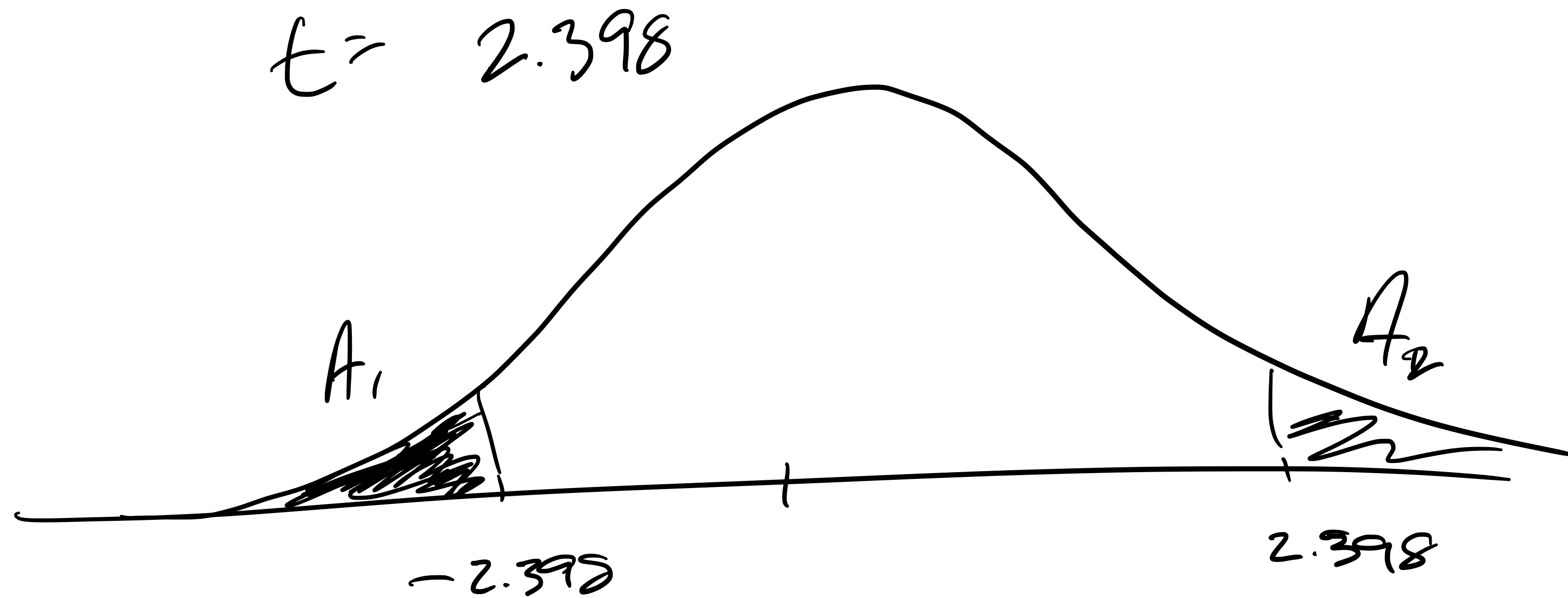| **Alternative Hypothesis** | **P-Value Level $\alpha$ Test** |
|:---:|:---:|
| $H_1 \ : \quad \theta > \theta_0$ | $P(T \geq t \mid H_0) \leq \alpha$ |
| $H_1 \ : \quad \theta < \theta_0$ | $P(T \leq t \mid H_0) \leq \alpha$ |
| $H_1 \ : \quad \theta \neq \theta_0$ | $2 \min \{ P(T \leq t \mid H_0), \ P(T \geq t \mid H_0) \} \leq \alpha$ |

# 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.

$$\bar{X}$$

$$S$$

$$\alpha = 0.1$$

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

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

$$= -2.398$$

# t-Test example (p-value method)

$t = \quad 2.398$



$A_1$                                                $A_2$

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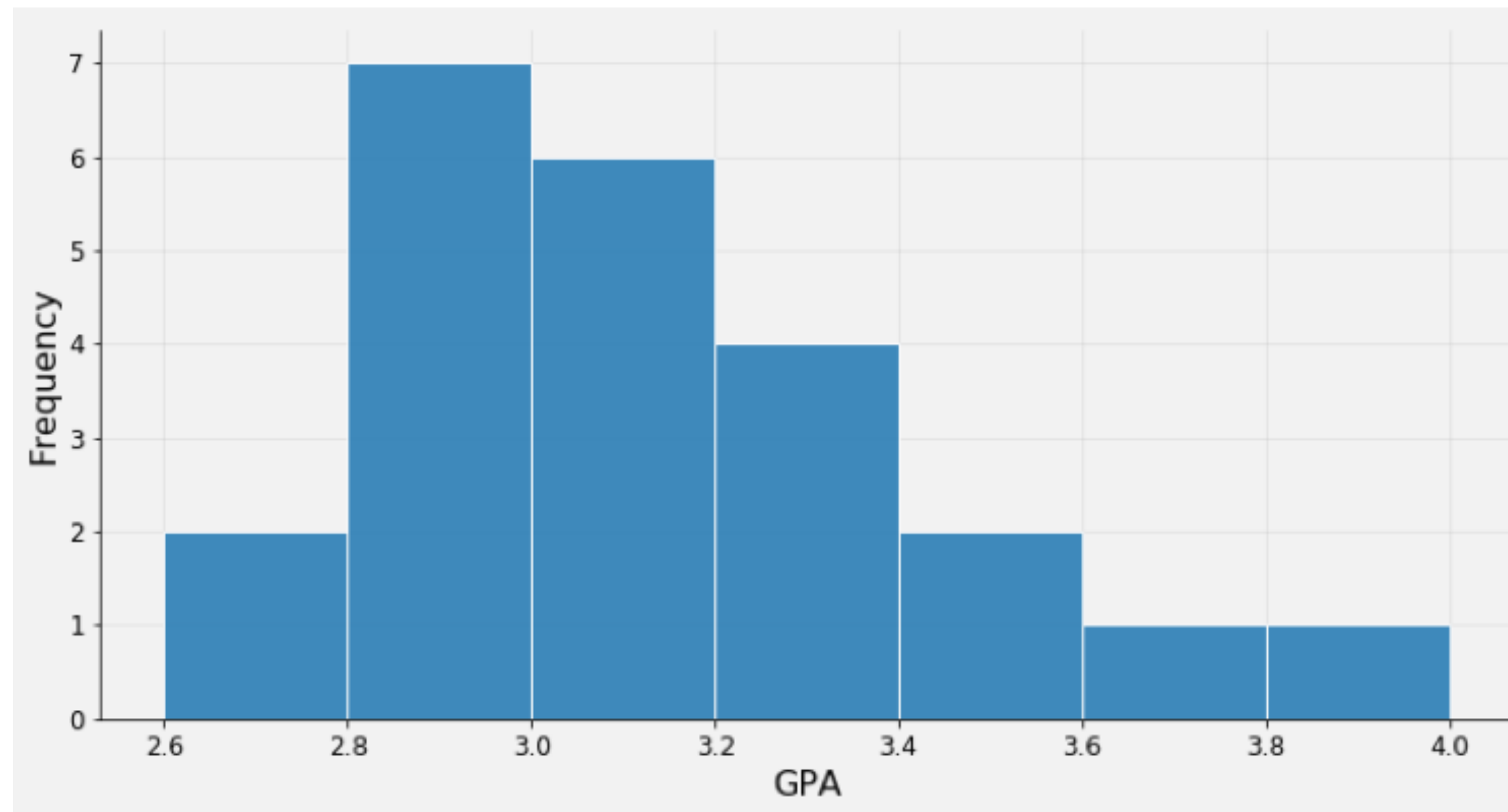
$$2 \times \text{stats.t.cdf}(\underset{\substack{\uparrow \\ \text{test} \\ \text{statistic}}}{-2.398}, \underset{\text{dof}}{22}) = \boxed{\underset{\text{p-value}}{0.0254}} < \underset{\alpha}{0.10}$$

# t-Test example (rejection region 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.
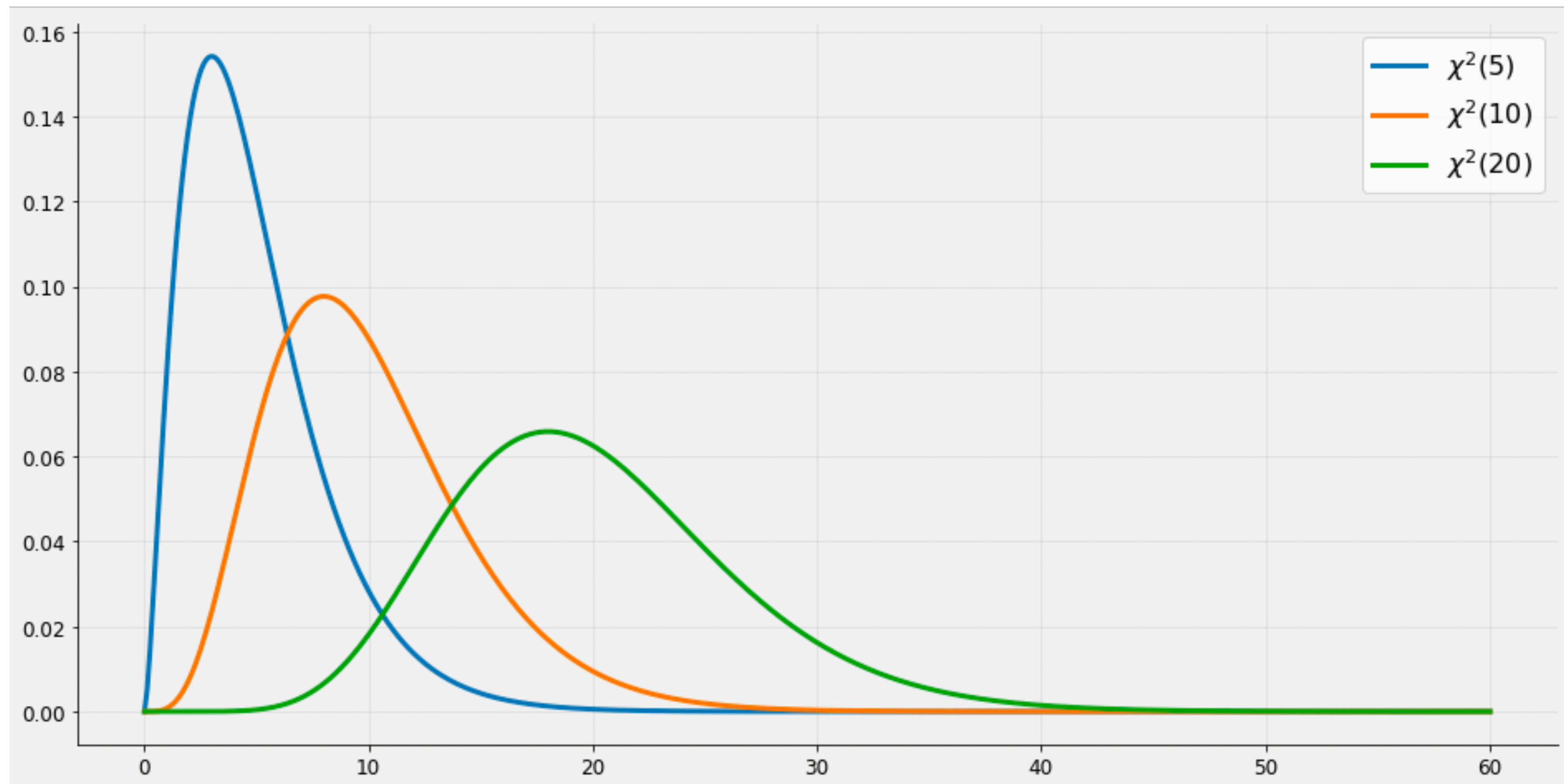
# t-Test example (rejection region method)

# Inference for *variances*

- After Spring Break, we'll talk about estimating confidence intervals for the variance of a population using something [wonderful] called **The Bootstrap**.

- But if your population is normally distributed, we have some [wonderful] theory which gives us a better confidence interval and works for both large and small sample sizes!

- **Question**: What does the sampling distribution of the variance look like when the population is **normally distributed**?

# The Chi-Squared Distribution

- The chi-squared distribution ( $\chi^2_\nu$ ) is also parameterized by degrees of freedom $\nu = n - 1$
- The pdfs of the family of $\chi^2_\nu$ distributions are gross, so lets just draw them!
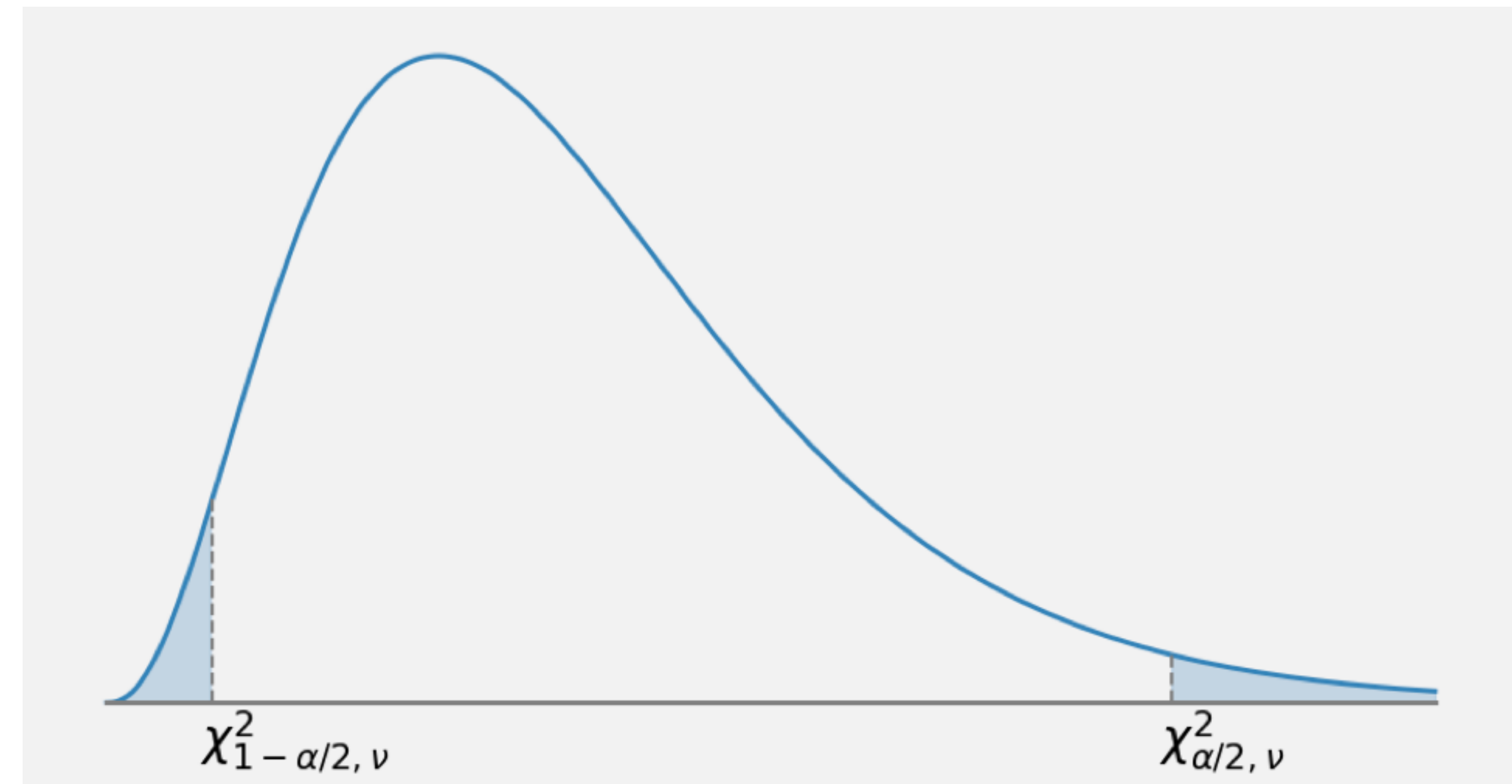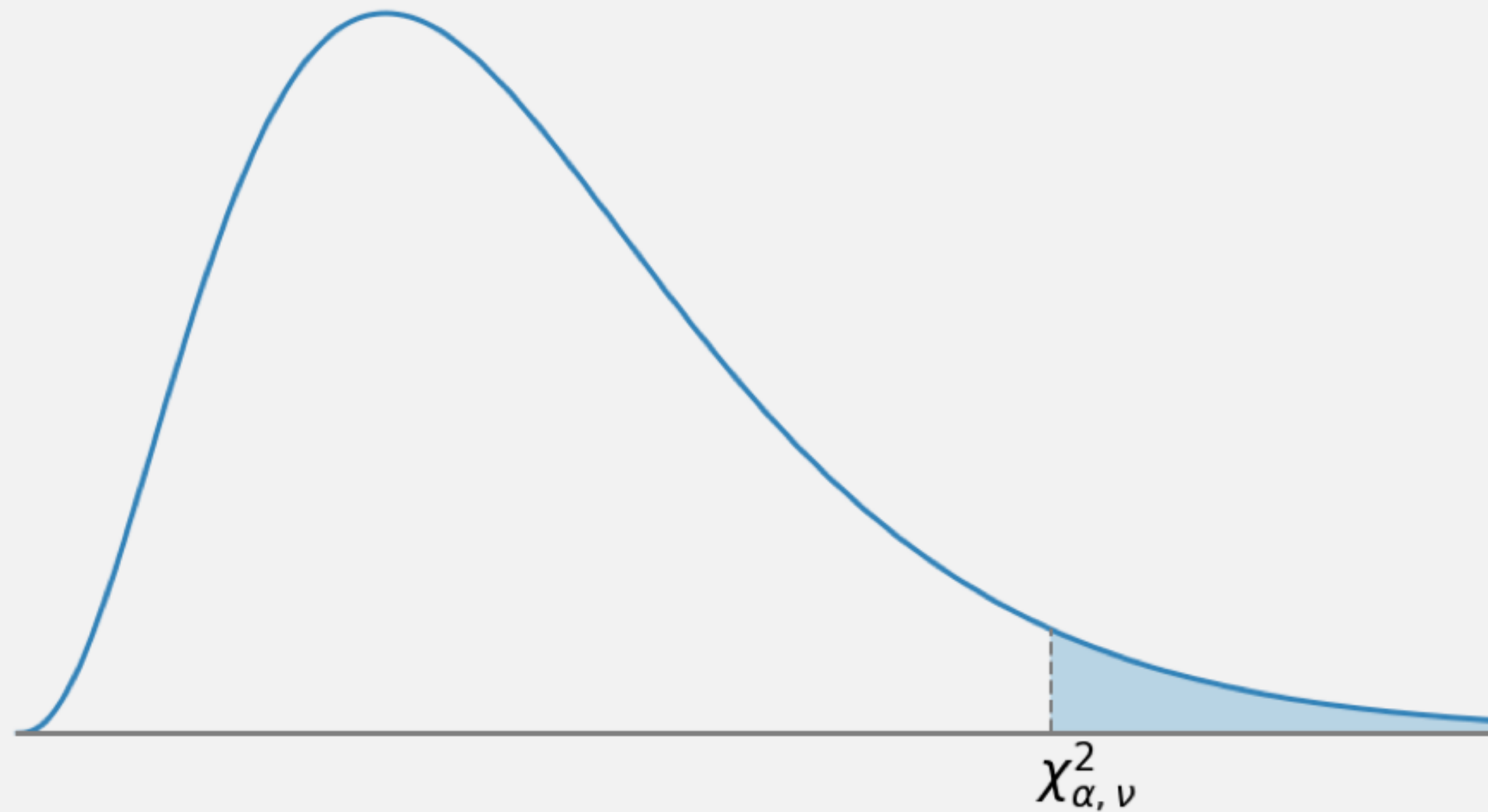
# A confidence interval for the variance

- Let $X_1, X_2, \ldots X_n$ be IID samples from a normal distribution with mean $\mu$ and standard deviation $\sigma$. Define the *sample variance* in the usual way as

- Then the random variable $(n-1)S^2/\sigma^2$ follows the distribution $\chi^2_{n-1}$.

- Then it follows that

# The Chi-Squared Dist is Non-Symmetric

- Because the distribution is non-symmetric, we need to use two different critical values.

# A confidence interval for the variance

- For a $100(1 - \alpha)\%$ confidence interval we choose the two critical values $X^2_{1-\alpha/2,n-1}$ and $X^2_{\alpha,n-1}$ which puts $\alpha/2$ probability in each tail. Then, with $100(1 - \alpha)\%$ confidence we can say that

# A confidence interval for the variance

- For a $100(1 - \alpha)\%$ confidence interval we choose the two critical values $X^2_{1-\alpha/2,n-1}$ and $X^2_{\alpha,n-1}$ which puts $\alpha/2$ probability in each tail. Then, with $100(1 - \alpha)\%$ confidence we can say that

$$\frac{(n - 1)S^2}{\chi^2_{\alpha/2,n-1}} < \sigma^2 < \frac{(n - 1)S^2}{\chi^2_{1-\alpha/2,n-1}}$$

**Question**: How can we use this to get a $100(1 - \alpha)\%$ confidence interval for the standard deviation?

- Example: A large candy manufacturer produces packages of candy targeted to weight 52g. The weight of the packages of candy is known to be normally distributed, but a QC engineer is concerned that the variation in the produced packages is larger than acceptable. In an attempt to estimate the variance she selects n=10 bags at random and weighs them. The sample yields a sample variance of 4.2g. Find a 95% confidence interval for the variance and a 95% confidence interval for the standard deviation.