



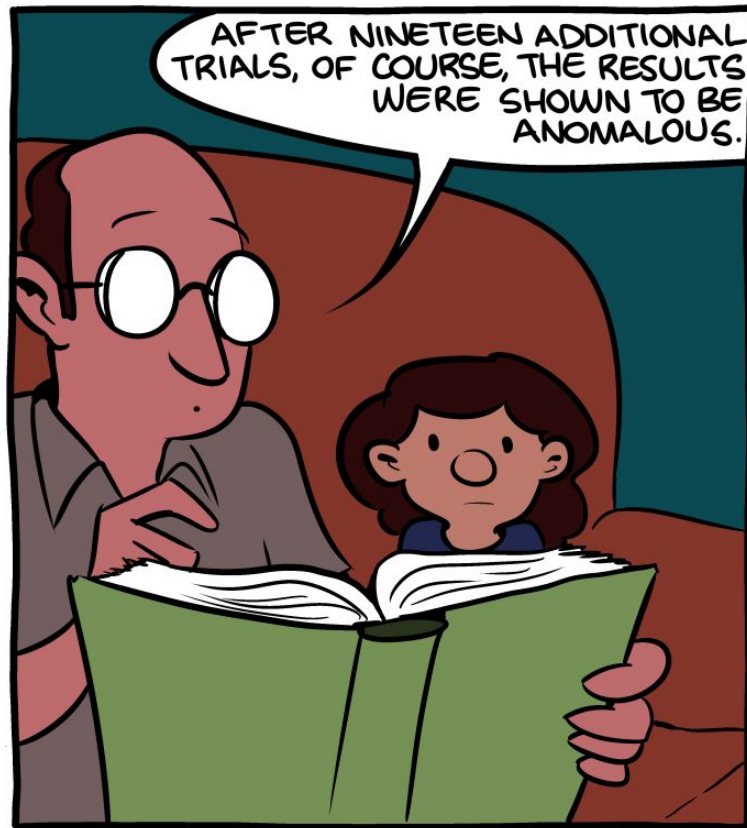
## Lecture 18: Statistical Inference with Small Samples



"The Tortoise And The Hare" is actually  
a fable about small sample sizes.

## Announcements and reminders

- HW 5 due **Friday 9 November at 5 PM**
- Check in on [Arkaive](#)!



"The Tortoise And The Hare" is actually a fable about small sample sizes.

## Previously, on CSCI 3022...

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

$\sigma$  is known:

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

$$E[\bar{x}] \text{ \& \; } \text{Var}(\bar{x})$$

$\sigma$  is unknown:

$$\frac{\bar{x} - \mu}{s / \sqrt{n}}$$
$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

## Previously, on CSCI 3022...

( $n \geq 30$ )

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

$\sigma$  is known:

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

$\sigma$  is unknown:

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

## Previously, on CSCI 3022...

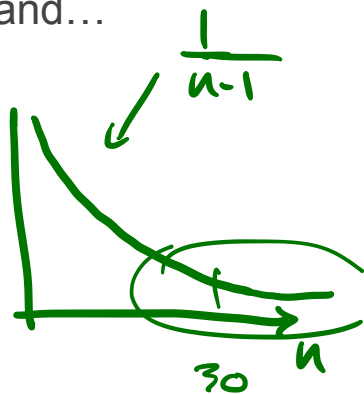
normality ←

$n < 30$

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

$\sigma$  is known:

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



$\sigma$  is unknown:

$\epsilon_i$  Normally distributed  
w/  $n < 30$

???







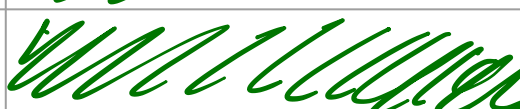

tempting...

$$\sqrt{\frac{1}{n-1} \sum (x_i - \bar{x})^2}$$

$$\frac{\bar{X} - \mu}{\underbrace{S}_{\leftarrow} / \sqrt{n}} \stackrel{?}{\sim} N(0, 1)$$

# 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

|  | $n \geq 30$  | $n < 30$  |
|--|--|---|
| Normal data, known $\sigma$ ✓          |  |  |
| Normal data, unknown $\sigma$          |  |  |
| <u>Non-normal</u> data, known $\sigma$ |  |  |
| Non-normal data, unknown $\sigma$      |  |  |

■ z-distribution  
(past)

■ t-distribution  
(today!)

■ Bootstrap (later)

## Small-sample tests for $\mu$

$$s^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2$$

When  $n$  is small, we can't invoke the Central Limit Theorem

- If we don't even know if the data are Normal, then we can **bootstrap**
- But that can be expensive (producing lots of replicates takes **time** and **memory**)

est. of  $\mu$  that  
we "bought" w/  
1 df

If we have **small  $n$**  and **some reason to think our data are (approximately) Normal**, then...

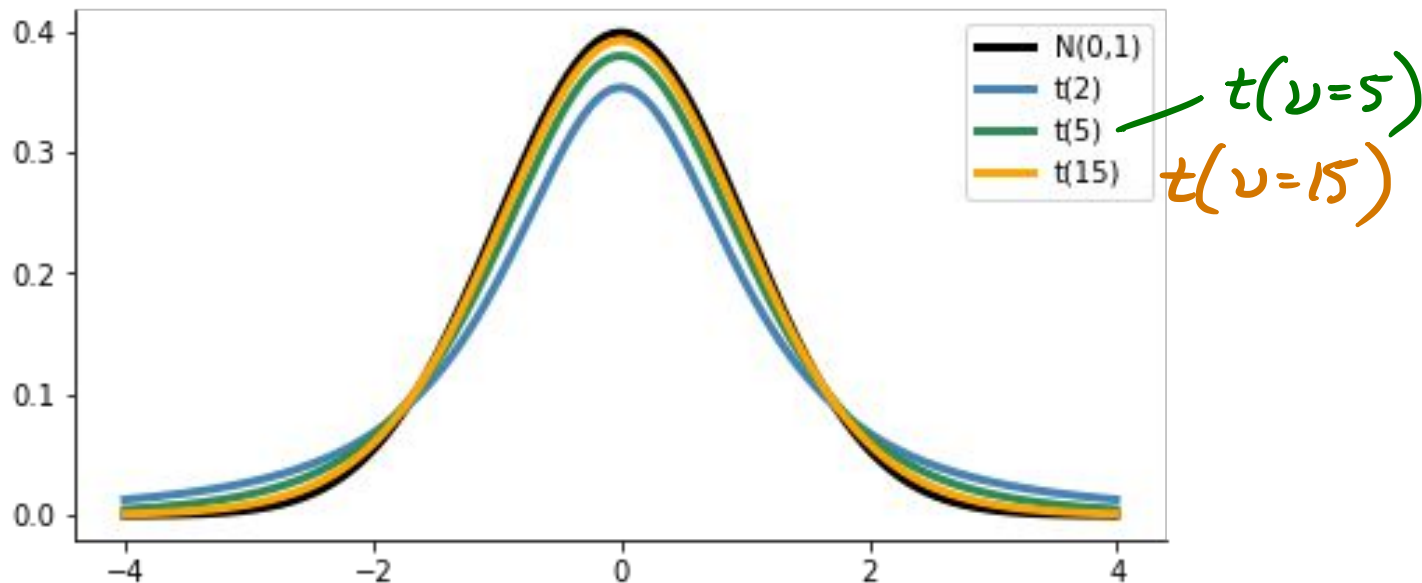
When  $\bar{X}$  is the sample mean of a random sample of size  $n$  from a normal distribution with mean  $\mu$ , the random variable

Test statistic: 
$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}} \sim t_v$$

follows a probability distribution called a **t-distribution** with parameter  $\nu = n-1$  degrees of freedom (df)

# The t-distribution

Here are some members of the family of t-distributions, and the standard normal  $N(0,1)$





# Properties of t-distributions

$n = \#$  data points (samples)



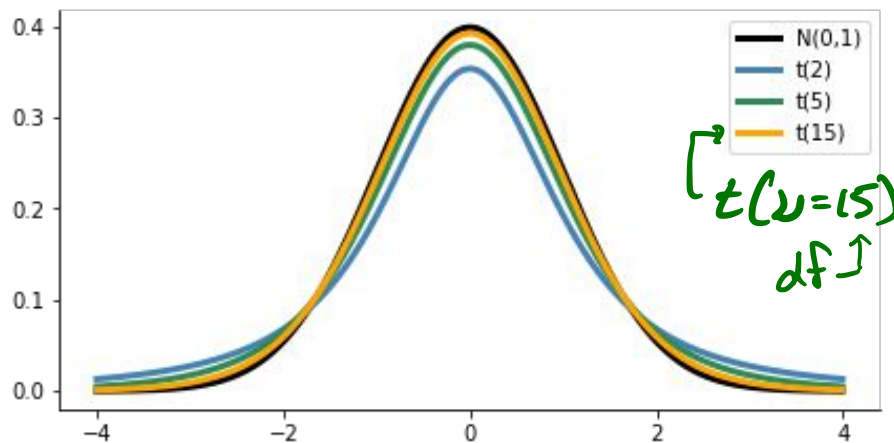
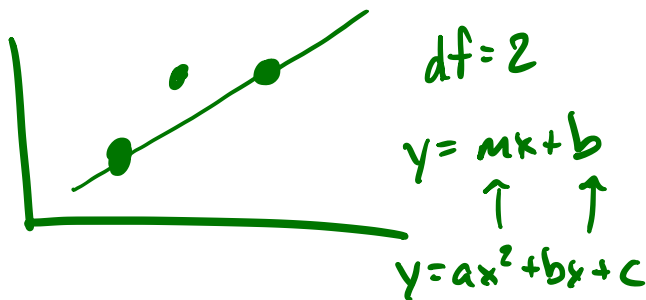
Let  $t_\nu$  denote the t-distribution with parameter  $\nu = n - 1$  df

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

\* 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  *$\lim_{\nu \rightarrow \infty} t_\nu = N(0,1)$*

• As  $\nu \rightarrow \infty$  the sequence of  $t_\nu$  curves approaches the standard normal curve

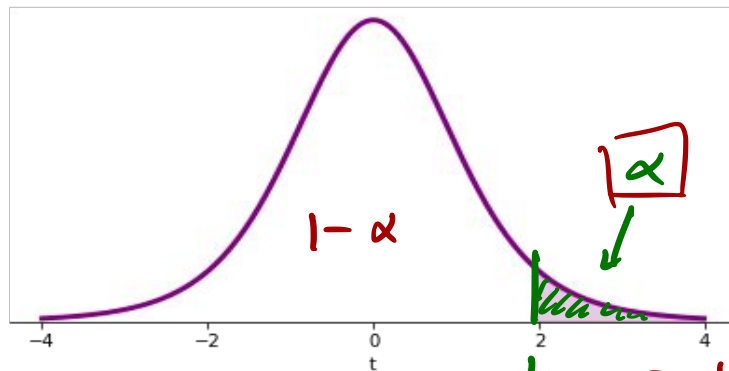


# The t-critical value

$\nu$

We can extend all of our inferential mechanics to the small-sample case by introducing the so-called t-critical value, which we denote as  $t_{\alpha, \nu}$

**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



2dof  
↓

$$t_{\alpha, \nu} = \text{stats.t.ppf}(1 - \alpha, \nu = n - 1)$$

**Example:**  $t_{0.05, 6}$  is the t-critical value that captures the upper-tail area of 0.05 (5%) under the t-curve with 6 degrees of freedom.

↑  $\nu = 6 = n - 1$

→ Sample size = 7

## The t-confidence interval for the mean

$$\left[ \bar{x} - z_{\alpha/2} \cdot \frac{s}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \cdot \frac{s}{\sqrt{n}} \right] \quad \text{--- } z\text{-dist.}$$

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

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

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

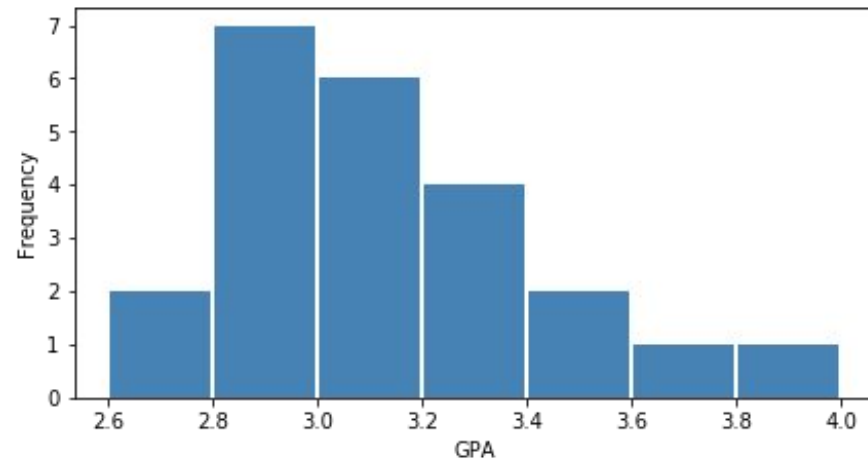
Or, more compactly:

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

# The t-confidence interval for the mean

$$\text{stats.t.ppf}(0.95, 22) = 1.717$$
$$[3.036, 3.256]$$

**Example:** S'pose the GPAs for 23 students have the histogram shown here. The sample mean of the data is 3.146 and the sample standard deviation is 0.308. Find a 90% CI for the mean GPA.



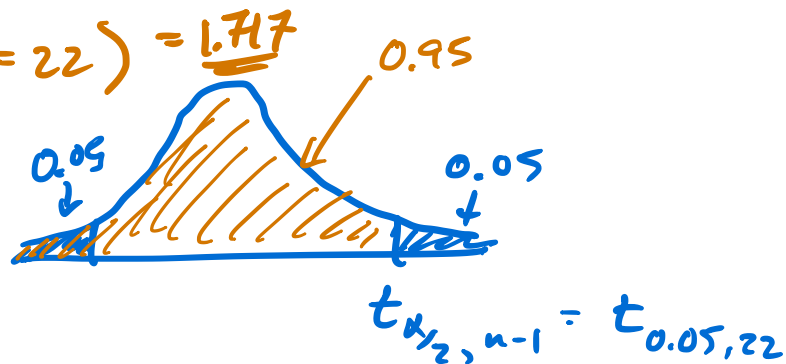
$\rightarrow \alpha = 0.1$

90% CI:  $\bar{X} \pm t_{0.05, 22} \cdot \frac{0.308}{\sqrt{23}}$

$\text{stats.t.ppf}(0.95, \text{ddof} = 22) = 1.717$

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

$$= [3.036, 3.256]$$



# The t-test, critical regions and p-values

$$H_0: \theta = \theta_0$$

Alternative hypothesis

Critical region level  $\alpha$  test

p-value level  $\alpha$  test

$$H_1: \theta > \theta_0$$

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

$$p\text{-value} = 1 - t.\text{cdf}(t, \text{ddof} = n-1)$$

$$H_1: \theta < \theta_0$$

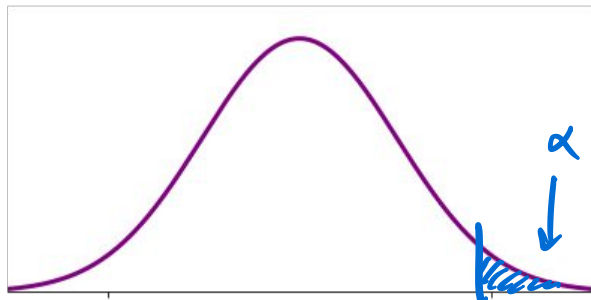
$$t \leq t_{1-\alpha, v}$$

$$p\text{-value} = t.\text{cdf}(t, \text{ddof} = n-1)$$

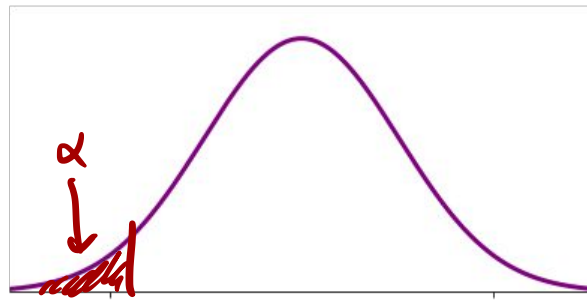
$$H_1: \theta \neq \theta_0$$

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

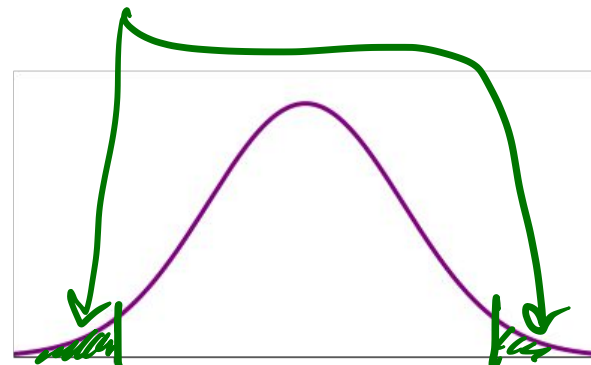
$$p\text{-value} = 2 \cdot t.\text{cdf}(-|t|, \text{ddof} = n-1)$$



$$t_{\alpha, v = n-1}$$



$$t_{1-\alpha, v}$$



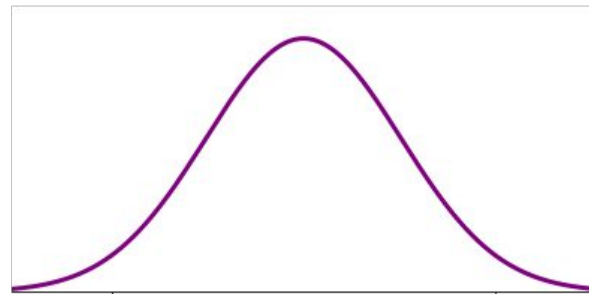
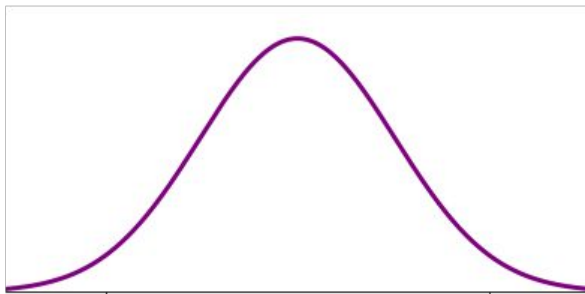
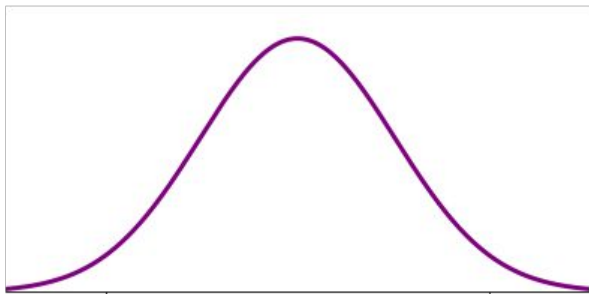
$$-|t|$$

$$|t|$$

$$\frac{\bar{X} - \theta_0}{s/\sqrt{n}} \rightarrow P(\tau > t)$$

# The t-test, critical regions and p-values

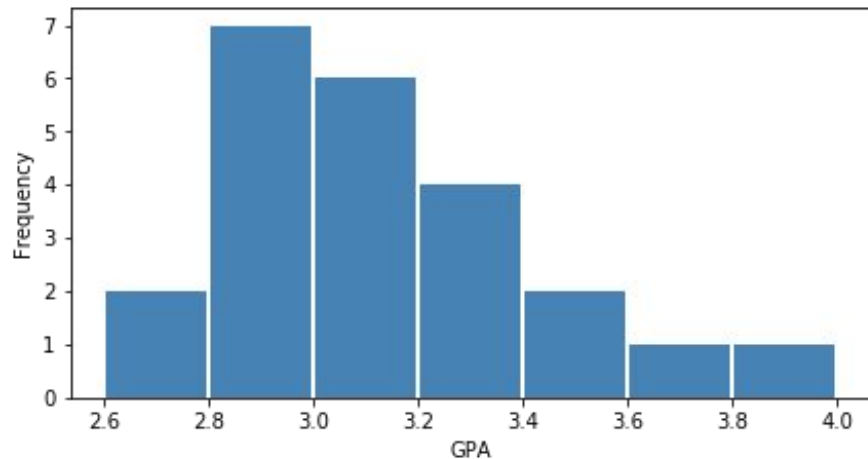
| Alternative hypothesis      | Critical region level $\alpha$ test                                  | p-value level $\alpha$ test  |
|-----------------------------|--|--|
| $H_1: \theta > \theta_0$    | $t \geq t_{\alpha, \nu}$   | $P(T \geq t \mid H_0) \leq \alpha$                                       |
| $H_1: \theta < \theta_0$    | $t \leq t_{\alpha, \nu}$   | $P(T \leq t \mid H_0) \leq \alpha$                                       |
| $H_1: \theta \neq \theta_0$ | $(t \geq t_{\alpha/2, \nu}) \text{ or } (t \leq -t_{\alpha/2, \nu})$ | $2 \cdot \min\{P(T \leq t \mid H_0), P(T \geq t \mid H_0)\} \leq \alpha$ |



## t-test for the mean, using p-values

$$\text{Stats.t.cdf}(-2.398, v=22) = 0.0127$$

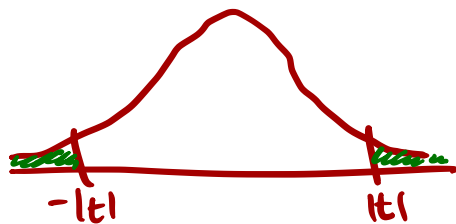
**Example:** S'pose the GPAs for 23 students have the histogram shown here. The sample mean of the data is 3.146 and the sample standard deviation is 0.308. Determine if there is sufficient evidence to conclude at the 0.10 (10%) significance level that the mean GPA is not equal to 3.30.



$$H_0: \mu = 3.30$$

$$H_1: \mu \neq 3.30$$

Two tailed



$$\alpha = 0.1$$

$$p\text{-value} = 2 \cdot t.\text{cdf}(-|t|, \text{dof}=22) = 2 \cdot 0.0127 = 0.0254201$$

Test statistic:

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

REJECT  $H_0$

# Inference for variances

---

We've talked about confidence intervals for the **mean** and for **proportions**

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

... if your population is **normally distributed**, it turns out we have some theory that gives us a **confidence interval** and works for both large *and* small samples!



## Inference for variances

---

We've talked about confidence intervals for the **mean** and for **proportions**

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

... if your population is **normally distributed**, it turns out we have some theory that gives us a **confidence interval** and works for both large *and* small samples!

## Inference for variances

---

We've talked about confidence intervals for the **mean** and for **proportions**

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

... if your population is **normally distributed**, it turns out we have some theory that gives us a **confidence interval** and works for both large ***and*** small samples!

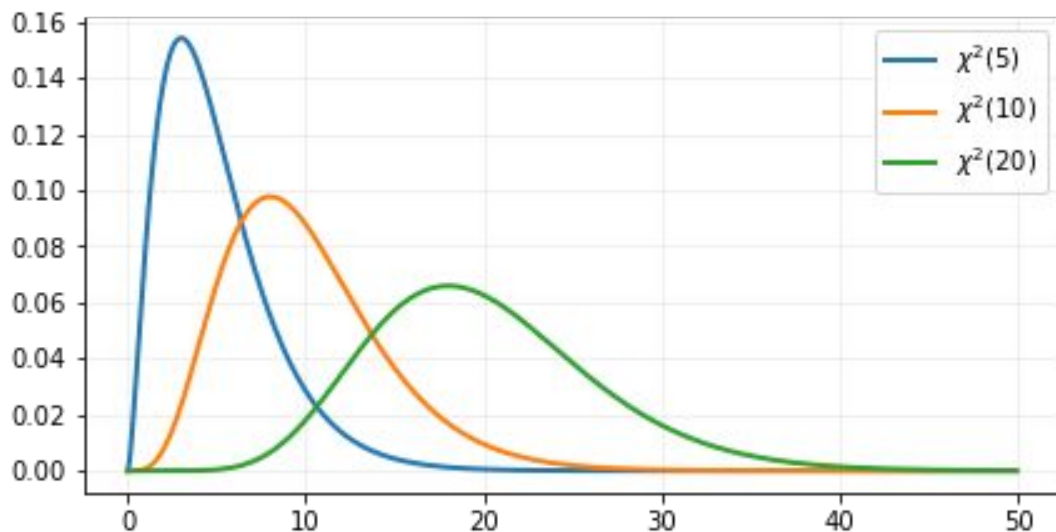
# The chi-squared distribution

( $\chi^2$  distribution)

$\chi^2$

The chi-squared ( $\chi^2$ ) distribution is also parameterized by degrees of freedom  $\nu = n-1$

The pdfs of the family  $\chi^2$  are pretty nasty, so let's just plot a few.



## A confidence interval for the variance

HERE 3p

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

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

Then the random variable  $\frac{(n-1) S^2}{\sigma^2}$  follows the distribution  $\chi_{n-1}^2$

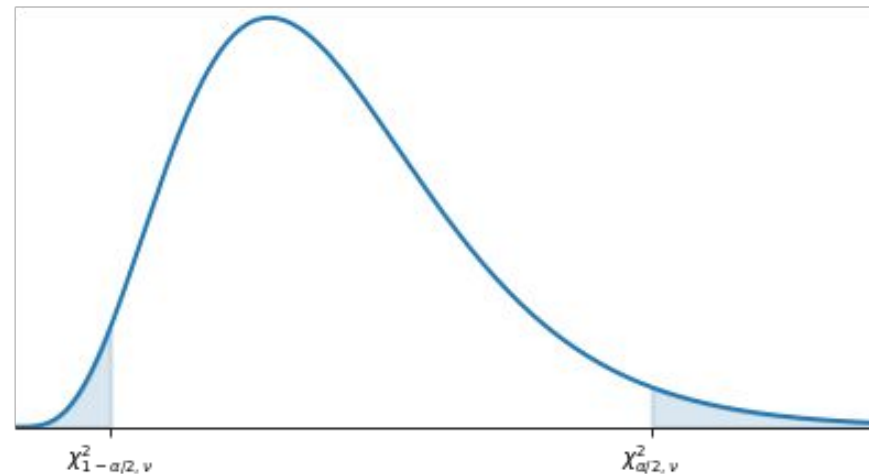
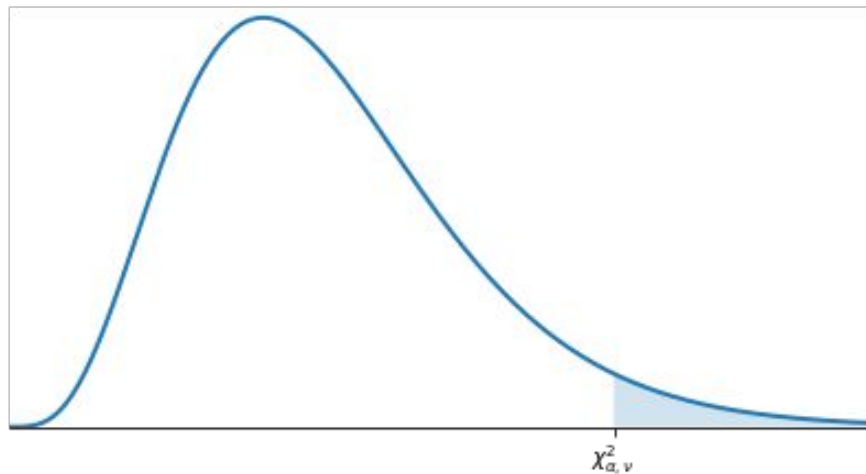
Then it follows that ...  $\hookrightarrow \frac{(n-1) S^2}{\sigma^2} \sim \chi_{n-1}^2$

$$\text{e.g. } P\left(\chi_{1-\alpha/2, n-1}^2 < \frac{(n-1) S^2}{\sigma^2} < \chi_{\alpha/2, n-1}^2\right) = 1 - \alpha$$

## A confidence interval for the variance

---

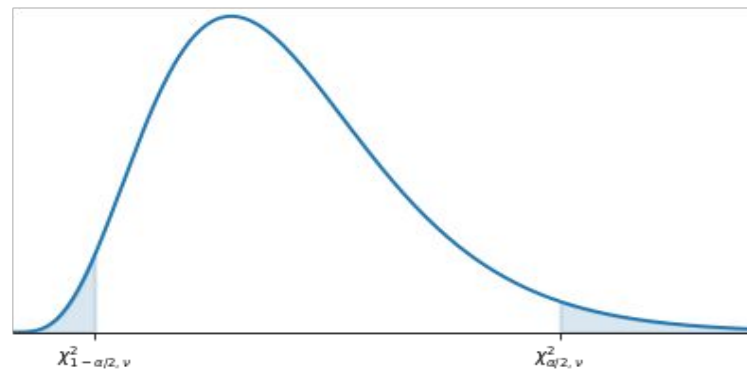
Because the  $\chi^2$  distribution is not symmetric, we need to use two different critical values



## A confidence interval for the variance

---

For a  $100 \cdot (1-\alpha)\%$  CI, we choose the **two** critical values  $\chi^2_{1-\alpha/2, n-1}$  and  $\chi^2_{\alpha/2, n-1}$ , which attributes  $\alpha/2$  probability to each the left and right tails. Then, with  $100 \cdot (1-\alpha)\%$  confidence we can say that



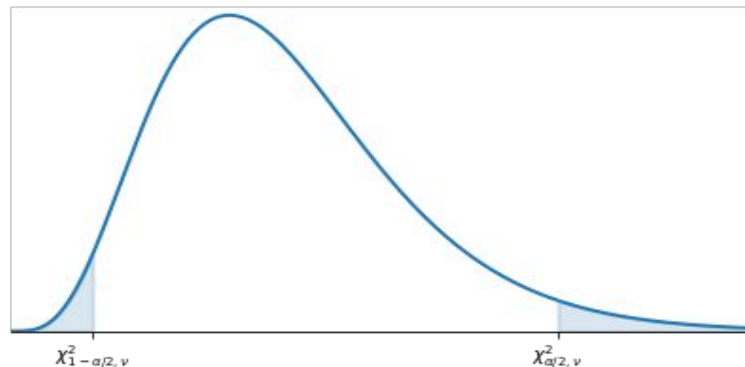
## A confidence interval for the variance

HERE

For a  $100 \cdot (1-\alpha)\%$  CI, we choose the **two** critical values  $\chi^2_{1-\alpha/2, n-1}$  and  $\chi^2_{\alpha/2, n-1}$ , which attributes  $\alpha/2$  probability to each the left and right tails. Then, with  $100 \cdot (1-\alpha)\%$  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:** What, then, is a  $100 \cdot (1-\alpha)\%$  CI for the SD?



## A confidence interval for the variance

---

**Example:** A large candy manufacturer produces packages of candy targeted to weigh 52 g. 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 he selects  $n=10$  bags at random and weighs them. The sample yields a sample variance of  $4.2 \text{ g}^2$ . Find a 95% CI for the variance, and a 95% CI for the SD.



## A confidence interval for the variance

---

**Example:** A large candy manufacturer produces packages of candy targeted to weigh 52 g. 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 he selects  $n=10$  bags at random and weighs them. The sample yields a sample variance of  $4.2 \text{ g}^2$ . Find a 95% CI for the variance, and a 95% CI for the SD.

$$\alpha = 0.05, \quad \alpha/2 = 0.025, \quad n = 10, \quad s^2 = 4.2$$

$$\chi_{1-\alpha/2, n-1}^2 = \chi_{0.975, 9}^2 = \text{stats.chi2.ppf}(0.025, 9) = 2.70$$

$$\chi_{\alpha/2, n-1}^2 = \chi_{0.025, 9}^2 = \text{stats.chi2.ppf}(0.975, 9) = 19.02$$

$$\frac{(10 - 1) \cdot 4.2}{19.02} < \sigma^2 < \frac{(10 - 1) \cdot 4.2}{2.70}$$
$$\Rightarrow 1.99 < \sigma^2 < 14.0$$

# What just happened?

- **Small samples** happened!
  - Learned what distributions (instead of standard normal) to use when our sample is too small for CLT to kick in
- **T-distributions** -- small sample CI/hypothesis testing for the **mean**
- **chi-squared distributions** -- small sample CI/hypothesis testing for the **variance**



"The Tortoise And The Hare" is actually a fable about small sample sizes.

