

# Inferences from Small Samples

# Outline

- Student's  $t$  distribution
- Small-sample inferences concerning a population mean
- Small-sample inferences for the difference between two population means: independent random samples
- Small-sample inferences for the difference between two means: a paired-difference test

# Student's $t$ Distribution

- Review of CLT results
  - If the population is normally distributed,  $\bar{x}$  and  $z$  follow a normal distribution regardless of sample size.
  - If the population is not normally distributed,  $\bar{x}$  and  $z$  follow a normal distribution if the sample size is large.
- When  $n$  is small ( $<30$ ) *and* the original population is not normally distributed, CLT does NOT guarantee that  $z$  will be normally distributed
  - The methods that we used for point and interval estimations and testing hypotheses no longer apply, e.g. the 95% confidence interval of  $\bar{x}$  is no longer  $\mu - 1.96SE < \bar{x} < \mu + 1.96SE$

# Student's $t$ Distribution

- The sampling distribution of  $\bar{x}$  and  $z$  can then be found by
  - Repeatedly drawing samples from the population, then computing and plotting the histogram of  $(\bar{x} - \mu)/(s/\sqrt{n})$
  - Deriving the actual distribution using the mathematical approach -> Student's  $t$  distribution.
- The distribution of statistic  $t = (\bar{x} - \mu)/(s/\sqrt{n})$  has the following characteristics
  - It has bell-shaped and symmetric around  $t$ , just like  $z$
  - It has more 'spread' than  $z$
  - It depends on the sample size. When  $n$  gets larger, the distribution of  $t$  becomes very similar to  $z$ .

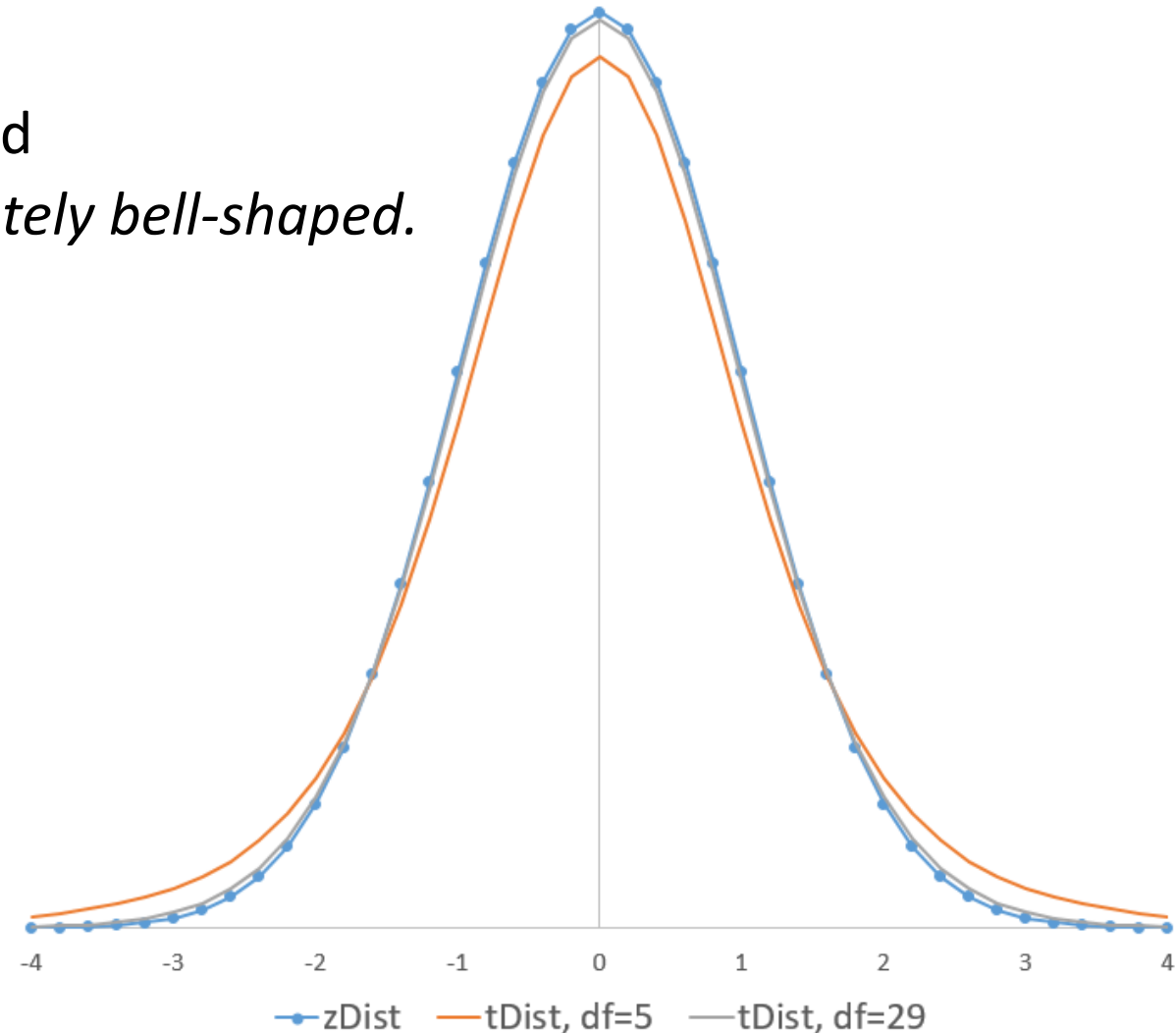
# Student's $t$ Distribution

- Conditions of Student's  $t$  distribution
  - Samples MUST be *randomly drawn* and
  - the population SHOULD be *approximately bell-shaped*.
- However,

“

Statisticians say that the  $t$  statistic is **robust**, meaning that the distribution of the statistic does not change significantly when the normality assumption is violated.

”



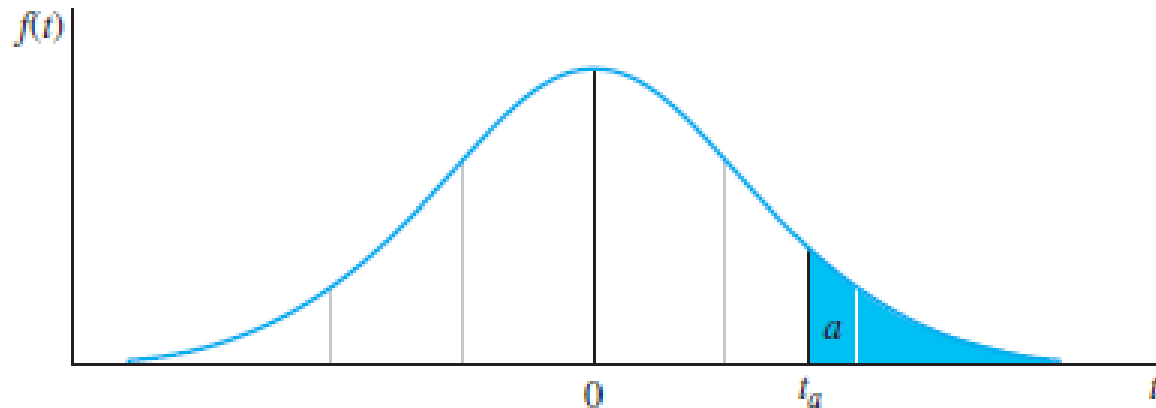
# Student's $t$ Distribution

**Example 1.** Calculate the probability of  $t > 2.015$  for  $df=5$ .

- In Excel,  $P(t > 2.015) = \text{TDIST}(2.015, 5, 1) = 0.05$ .

**Example 2.** Calculate the  $t$  value larger than 1% of all values of  $t$  for  $df=9$ .

- In Excel,  $t = \text{T.INV}(0.01, 9) = -2.821$



$df$	$t_{.100}$	$t_{.050}$	$t_{.025}$	$t_{.010}$	$t_{.005}$	$df$
1	3.078	6.314	12.706	31.821	63.657	1
2	1.886	2.920	4.303	6.965	9.925	2
3	1.638	2.353	3.182	4.541	5.841	3
4	1.533	2.132	2.776	3.747	4.604	4
5	1.476	2.015	2.571	3.365	4.032	5
6	1.440	1.943	2.447	3.143	3.707	6
7	1.415	1.895	2.365	2.998	3.499	7
8	1.397	1.860	2.306	2.896	3.355	8
9	1.383	1.833	2.262	2.821	3.250	9
.	.	.	.	.	.	.
.	.	.	.	.	.	.
.	.	.	.	.	.	.
26	1.315	1.706	2.056	2.479	2.779	26
27	1.314	1.703	2.052	2.473	2.771	27
28	1.313	1.701	2.048	2.467	2.763	28
29	1.311	1.699	2.045	2.462	2.756	29
inf.	1.282	1.645	1.960	2.326	2.576	inf.

# Small-sample Inferences for a Population Mean

Apply the same procedure as in estimation and hypothesis testing for large samples

- $(1 - \alpha)\%$  confidence interval for  $\mu$  is  $\bar{x} \pm t_{\alpha/2}s/\sqrt{n}$

- Hypothesis testing:

(1) Null hypothesis

$$H_0: \mu = \mu_0$$

(2) Alternative hypothesis

$$H_a: \mu > \mu_0, \text{ or } H_a: \mu < \mu_0 \text{ (one-tailed test),}$$

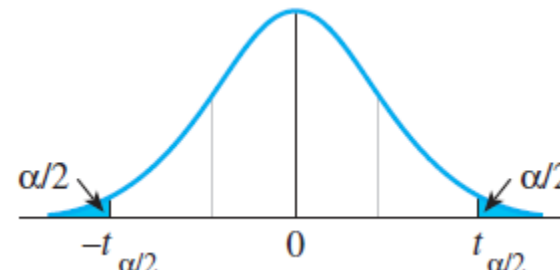
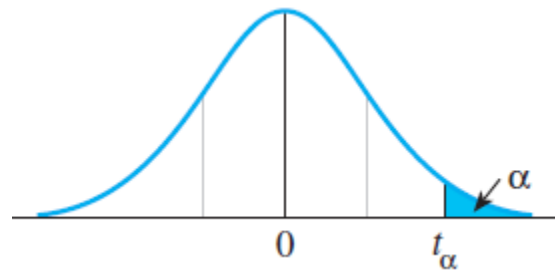
$$H_a: \mu \neq \mu_0 \text{ (two-tailed test)}$$

(3) Test statistic

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

(4) Rejection region (Note that the critical values of  $t$  is based on  $(n-1)$  degrees of freedom)

$$t > t_\alpha \text{ (for } H_a: \mu > \mu_0)$$
$$t < -t_\alpha \text{ (for } H_a: \mu < \mu_0)$$



$$t > t_{\alpha/2} \text{ or } t < -t_{\alpha/2}$$
$$\text{(for } H_a: \mu \neq \mu_0 \text{)}$$

OR pValue  $< \alpha$

# Small-sample Inferences for a Population Mean

**Example 3.** A paint manufacturer claimed that a can of 3.78l of their paint can cover 37.2 m<sup>2</sup> of wall area. In order to test this claim, 10 random cans were used to paint on 10 identical areas using the same kind of equipment. The actual area (in m<sup>2</sup>) covered by each of the 10 cans are as below

28.8	28.9	38.3	34.2	41.5
34.9	28.1	38.1	33.9	32.5

Does the test present sufficient evidence to support the manufacturer claim? Use  $\alpha = .05$ . Calculate the 95% confidence interval of the coverable area based on the test data.

- (1) Null hypothesis
- (2) Alternative hypothesis
- (3) Test statistic
- (4) Rejection region
- (5) Conclusion



# Small-sample inferences for the difference between 2 population means

Apply the same procedure as in estimation and hypothesis testing for large samples

- $(1 - \alpha)\%$  confidence interval for  $\mu_1 - \mu_2$  is  $(\bar{x}_1 - \bar{x}_2) \pm t_{\alpha/2} s \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$

- Hypothesis testing

(1) Null hypothesis  $H_0$ :

$$\mu_1 - \mu_2 = D_0$$

(2) Alternative hypothesis

$$H_a: \mu_1 - \mu_2 > D_0, \text{ or } H_a: \mu_1 - \mu_2 < D_0 \text{ (one-tailed test)}$$

$$H_a: \mu_1 - \mu_2 \neq D_0 \text{ (two-tailed test)}$$

(3) Test statistic

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{s \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}, \text{ where } s^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

# Small-sample inferences for the difference between 2 population means

(4) Rejection region (Note that the critical values of  $t$  is based on  $(n_1 + n_1 - 2)$  degrees of freedom)

$$t > t_{\alpha} \text{ (for } H_a: \mu_1 - \mu_2 > D_0)$$

$$t > t_{\alpha/2} \text{ or } t < -t_{\alpha/2} \text{ (} H_a: \mu_1 - \mu_2 \neq D_0)$$

$$t < -t_{\alpha} \text{ (for } H_a: \mu_1 - \mu_2 < D_0)$$

Or  $\text{pValue} < \alpha$

- Assumptions:
  - Samples must be *randomly selected*
  - Samples must be *independent*
  - Population variances must be *equal* or *nearly equal*.

# Small-sample inferences for the difference between 2 population means

**Example 4.** The time required by two swimmers to complete each of 10 trials of 100m freestyle swimming were recorded as below.

Swimmer 1	59.62	59.48	59.65	59.5	60.01	59.74	59.43	59.72	59.63	59.68
Swimmer 2	59.81	59.32	59.76	59.64	59.86	59.41	59.63	59.5	59.83	59.51

Do the data provide sufficient evidence to conclude that one swimmer is faster than the other?

# Small-sample inferences for the difference between 2 population means

- In cases where the two variances *are significantly different*, e.g. *Larger  $s^2$  / Smaller  $s^2 > 3$* , the above formulae for hypothesis testing of 2 population means need revisions, as below.

$$\text{Test statistic: } t = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \quad \text{Degree of freedom} \approx \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\frac{(s_1^2/n_1)^2}{n_1-1} + \frac{(s_2^2/n_2)^2}{n_2-1}}$$

**Example 5.** The number of raisins in 14 random miniboxes of Sunmaid® and in 14 random miniboxes of a generic brand were counted and presented below.

Generic Brand	25	26	25	28	26	28	28	27	26	27	24	25	26	26
Sunmaid®	25	29	24	24	28	24	28	22	25	28	30	27	28	24

Is there enough evidence to conclude that there is a significant difference between the average number of raisins in miniboxes of Sunmaid® and of the generic brand?

# Small-sample inferences for the difference between 2 population means – Paired-difference test

**Example 6.** Ten randomly selected drivers were shown a prohibitive sign of ‘No Left Turn’, and a permissive sign of ‘Left Turn Only’ during a driver reaction test. Their response time (in ms) to each of the signs were recorded and are presented below.

Driver	1	2	3	4	5	6	7	8	9	10
No left turn	824	866	841	770	829	764	857	831	846	759
Left turn only	702	725	744	663	792	708	747	685	742	610

Is there enough evidence to conclude that there is a difference between the response time to ‘No left turn’ sign and to ‘Left turn only’ sign?

**But, is this test different to the other tests presented above?**

# Small-sample inferences for the difference between 2 population means – Paired-difference test

- Paired-difference tests help reduce the effect of potential large variability among experimental units.
- The two samples are no longer independent.
- Use the same procedure for hypothesis testing and estimation of population mean

(1) Null hypothesis

$$H_0: \mu_D = 0 \text{ where } \mu_D = \mu_1 - \mu_2$$

(2) Alternative hypothesis

$$H_a: \mu_D > 0, \text{ or } H_a: \mu_D < 0 \quad (\text{one-tailed test})$$

$$H_a: \mu_D \neq 0 \quad (\text{two-tailed test})$$

(3) Test statistic  $t = \frac{\bar{d}}{s_d/\sqrt{n}}$

where  $n$  = number of paired differences

$\bar{d}$  = mean of the sample differences

$s_d$  = standard deviation of the sample differences

(4) Rejection region

$$t > t_\alpha \text{ (for } H_a: \mu_D > 0)$$

$$t < t_{\alpha/2} \text{ or } t > -t_{\alpha/2} \text{ (for } H_a: \mu_D \neq 0)$$

$$t < -t_\alpha \text{ (for } H_a: \mu_D < 0)$$

OR  $p\text{Value} < \alpha$  (Note the degree of freedom is  $n-1$ )

# Small-sample inferences for the difference between 2 population means – Paired-difference test

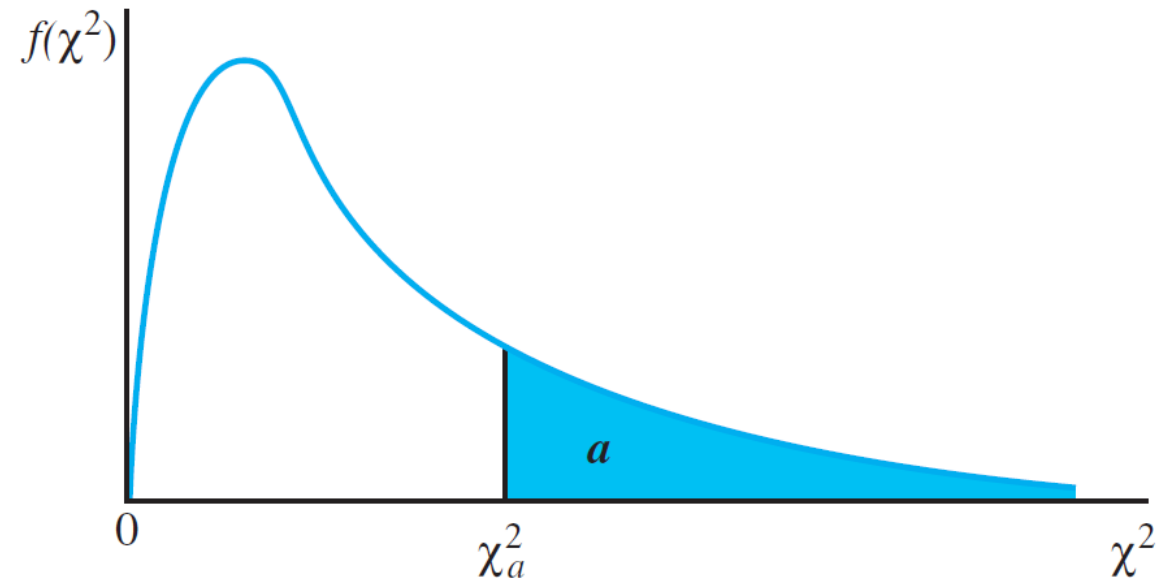
## Example 6 (cont.)

Driver	1	2	3	4	5	6	7	8	9	10
No left turn	824	866	841	770	829	764	857	831	846	759
Left turn only	702	725	744	663	792	708	747	685	742	610
Sample difference	122	141	97	107	37	56	110	146	104	149

- (1) Null hypothesis
- (2) Alternative hypothesis
- (3) Test statistic
- (4) Rejection region
- (5) Conclusion

# Inferences concerning a population variance

- The sampling distribution of  $s^2 = \frac{\sum(x_i - \bar{x})^2}{n-1}$  has the following behaviours:
  - Has its mean equal to  $\sigma^2$ , i.e. is an un-biased estimator of  $\sigma^2$
  - Has a lower bound of 0 (variance cannot be negative)
  - Is nonsymmetric
  - Has its shape depending on n
- Standardised variance  $\chi^2 = \frac{(n-1)s^2}{\sigma^2}$
- Chi-square probability distribution

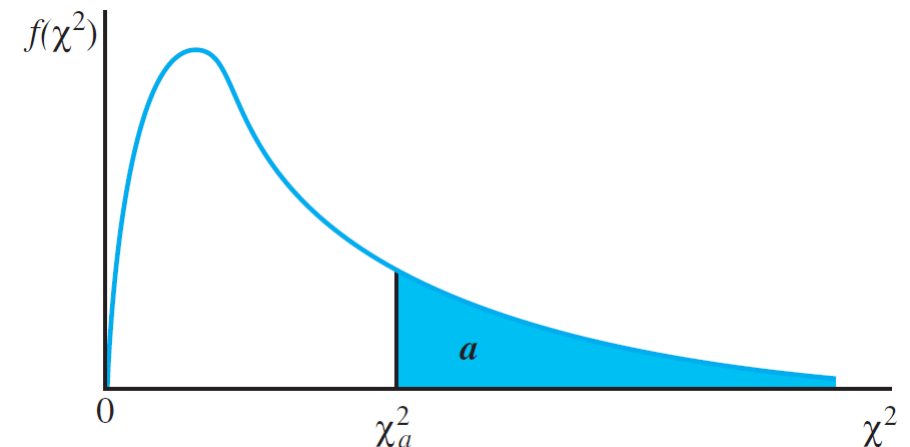




# Inferences concerning a population variance

$df$	$\chi^2_{.995}$	...	$\chi^2_{.950}$	$\chi^2_{.900}$	$\chi^2_{.100}$	$\chi^2_{.050}$	...	$\chi^2_{.005}$	$df$
1	.0000393		.0039321	.0157908	2.70554	3.84146		7.87944	1
2	.0100251		.102587	.210720	4.60517	5.99147		10.5966	2
3	.0717212		.351846	.584375	6.25139	7.81473		12.8381	3
4	.206990		.710721	1.063623	7.77944	9.48773		14.8602	4
5	.411740		1.145476	1.610310	9.23635	11.0705		16.7496	5
6	.0675727		1.63539	2.204130	10.6446	12.5916		18.5476	6
.	.		.	.	.	.		.	.
.	.		.	.	.	.		.	.
.	.		.	.	.	.		.	.
15	4.60094		7.26094	8.54675	22.3072	24.9958		32.8013	15
16	5.14224		7.96164	9.31223	23.5418	26.2962		34.2672	16
17	5.69724		8.67176	10.0852	24.7690	27.5871		35.7185	17
18	6.26481		9.39046	10.8649	25.9894	28.8693		37.1564	18
19	6.84398		10.1170	11.6509	27.2036	30.1435		38.5822	19

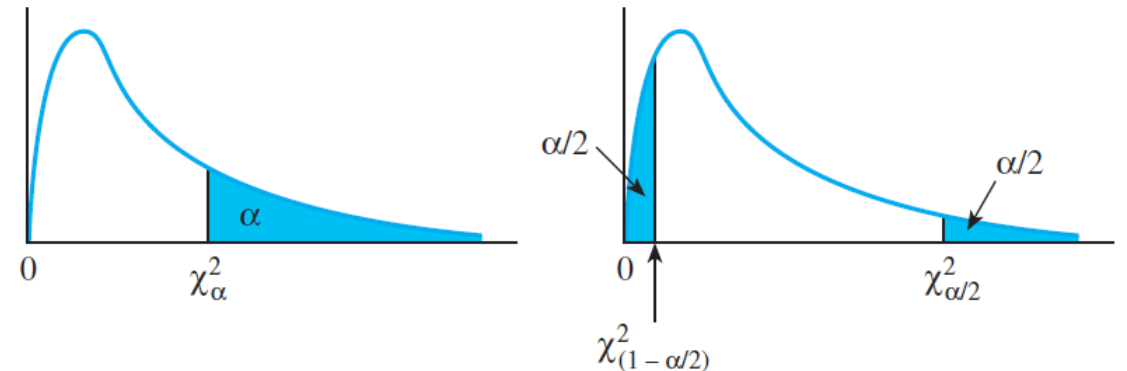
- In Excel,
  - CHIDIST(24.9958,15) = .05
  - CHIINV(0.05,15) = 24.99579



# Inferences concerning a population variance

- Hypothesis testing of a population variance
  - Null hypothesis  $\mathbf{H_0: \sigma^2 = \sigma_0^2}$
  - Alternative hypothesis  $\mathbf{H_a: \sigma^2 > \sigma_0^2}$  or  $\mathbf{H_a: \sigma^2 < \sigma_0^2}$  (One-tailed test)  
 $\mathbf{H_0: \sigma^2 \neq \sigma_0^2}$  (Two-tailed test)
  - Test statistic  $\chi^2 = \frac{(n-1)s^2}{\sigma_0^2}$
  - Rejection region: Reject  $\mathbf{H_0}$  when
$$\chi^2 > \chi_{\alpha}^2 \text{ (or } \chi^2 > \chi_{(1-\alpha)}^2 \text{)} \quad \text{(One-tailed test)}$$
$$\chi^2 > \chi_{\alpha/2}^2 \text{ or } \chi^2 < \chi_{(1-\frac{\alpha}{2})}^2 \quad \text{(Two-tailed test)}$$

Or when p-value  $\leq \alpha$
  - Conclusions



# Inferences concerning a population variance

**Example 7.** The user manual of a measuring instrument claimed that its variability was measured by a standard deviation of 2. In an experiment, the experimenter recorded reading values of 4.1, 5.2 and 10.2. Does the data confirm or disprove the manual's claim of the instrument's variability?

- (a) Test the appropriate hypothesis with  $\alpha = .1$
- (b) Calculate the 90% confidence interval of the true variance.