

The next example has random sampling from two normal populations, and the interest is in knowing whether there is a statistically significant difference between the population means of the two normal populations.

**Example 3.11** The tomato yields (in pounds) for Fertilizer  $X$  and Fertilizer  $Y$  are given in Table 3.2. Find a 90% confidence interval for the difference between the population means for the two fertilizers.

Fertilizer $X$ yields	29.9	11.4	25.3	16.5	21.2	
Fertilizer $Y$ yields	26.6	23.7	28.5	14.2	17.9	24.3

Table 3.2: Tomato yields for Fertilizer  $X$  and Fertilizer  $Y$ .

Fertilizer  $Y$  is a proposed fertilizer and the interest is whether it outperforms Fertilizer  $X$  in terms of expected yield. Some sample statistics from this data set are

$$\bar{x} = 20.9 \quad \bar{y} = 22.5 \quad s_X = 7.25 \quad s_Y = 5.43.$$

There is an increased yield of 1.6 pounds by using Fertilizer  $Y$  instead of Fertilizer  $X$ . But there is also significant spread to the data values as evidenced by the two standard deviations. So is the 1.6 pound increase in the yield a matter of random sampling variability or is it due to superior performance by Fertilizer  $Y$ ?

There are a host of problems associated with constructing the 90% confidence interval for  $\mu_X - \mu_Y$ . There are only  $n = 5$  observed yields from Fertilizer  $X$  and  $m = 6$  observed yields from Fertilizer  $Y$ . These sample sizes are too small to plot a histogram to determine whether the normality assumption is justified. On the other hand, previous experience with tomato yields might allow one to conclude that the assumption of normality is reasonable for tomato yields. For now, let's assume that the normality assumption is justified based on previous experiments. The appropriate formula for a confidence interval for the difference between the population means from Table 3.1 is the confidence interval associated with Theorem 1.7. This confidence interval not only requires that the random samples come from normally distributed populations—it also requires that the population variances are equal. Although the sample standard deviations,  $s_X = 7.25$  and  $s_Y = 5.43$ , are nearly equal, we would like some objective measure to indicate that they could indeed be considered equal for constructing the confidence interval. One way to proceed is to calculate a confidence interval for the ratio of the variances using the entry associated with Theorem 1.8 from Table 3.1. So an exact two-sided 90% confidence interval for the ratio of the population variances is

$$\frac{s_X^2}{s_Y^2 F_{n-1, m-1, \alpha/2}} < \frac{\sigma_X^2}{\sigma_Y^2} < \frac{s_X^2 F_{m-1, n-1, \alpha/2}}{s_Y^2}$$

or

$$\frac{7.25^2}{5.43^2 F_{4, 5, 0.05}} < \frac{\sigma_X^2}{\sigma_Y^2} < \frac{7.25^2 F_{5, 4, 0.05}}{5.43^2}$$

or

$$0.343 < \frac{\sigma_X^2}{\sigma_Y^2} < 11.1.$$

This confidence interval can be calculated with the following four R statements.

```
x = c(29.9, 11.4, 25.3, 16.5, 21.2)
y = c(26.6, 23.7, 28.5, 14.2, 17.9, 24.3)
l = var(x) / (var(y) * qf(0.95, 4, 5))
u = var(x) * qf(0.95, 5, 4) / var(y)
```

Since this confidence interval covers 1 by a wide margin, it is reasonable to conclude that the differences between the two sample standard deviations can be attributed to random sampling variability rather than a difference in the population variances. We can now proceed with the calculation of the confidence interval for the differences in the yields. The pooled sample variance is

$$s_p^2 = \frac{(n-1)s_X^2 + (m-1)s_Y^2}{n+m-2} = \frac{4 \cdot 7.25^2 + 5 \cdot 5.43^2}{5+6-2} = 39.7.$$

The associated pooled sample standard deviation is  $s_p = 6.30$ . So assuming normality and equal population variances in the two populations, the exact two-sided 90% confidence interval for  $\mu_X - \mu_Y$  is

$$\bar{x} - \bar{y} - t_{n+m-2, \alpha/2} s_p \sqrt{\frac{1}{n} + \frac{1}{m}} < \mu_X - \mu_Y < \bar{x} - \bar{y} + t_{n+m-2, \alpha/2} s_p \sqrt{\frac{1}{n} + \frac{1}{m}}$$

or

$$20.9 - 22.5 - t_{9, 0.05} 6.30 \sqrt{\frac{1}{5} + \frac{1}{6}} < \mu_X - \mu_Y < 20.9 - 22.5 + t_{9, 0.05} 6.30 \sqrt{\frac{1}{5} + \frac{1}{6}}$$

or

$$-8.67 < \mu_X - \mu_Y < 5.32.$$

This confidence interval can be calculated with the following R code.

```
x = c(29.9, 11.4, 25.3, 16.5, 21.2)
y = c(26.6, 23.7, 28.5, 14.2, 17.9, 24.3)
n = length(x)
m = length(y)
s = sqrt(((n - 1) * var(x) + (m - 1) * var(y)) / (n + m - 2))
l = mean(x) - mean(y) - qt(0.95, n + m - 2) * s * sqrt(1 / n + 1 / m)
u = mean(x) - mean(y) + qt(0.95, n + m - 2) * s * sqrt(1 / n + 1 / m)
```

Since confidence intervals of this type are calculated routinely by statisticians, some keystrokes can be saved by using the built-in function `t.test` with the specified stated coverage included in the `conf.level` parameter, and the `var.equal` parameter set to `TRUE`, as shown below.

```
x = c(29.9, 11.4, 25.3, 16.5, 21.2)
y = c(26.6, 23.7, 28.5, 14.2, 17.9, 24.3)
t.test(x, y, conf.level = 0.90, var.equal = TRUE)
```

Since this confidence interval includes 0, it can be concluded that there is not enough statistical evidence here to conclude that Fertilizer *Y* is better than Fertilizer *X*. Larger sample sizes for the two populations might be needed to conclude that the difference between the population means is nonzero.

This concludes the discussion of exact confidence intervals. The next section considers the case when  $P(L < \theta < U)$  does not equal  $1 - \alpha$ . These confidence intervals are known as approximate confidence intervals.