



## CHAPTER 15

- 15.1 The Wilcoxon Rank Sum Test
- 15.2 The Wilcoxon Signed Rank Test
- 15.3 The Kruskal-Wallis Test

# Nonparametric Tests

## Introduction

The most commonly used methods for inference about the means of quantitative response variables assume that the variables in question have Normal distributions in the population or populations from which we draw our data. In practice, of course, no distribution is exactly Normal. Fortunately, our usual methods for inference about population means (the one-sample and two-sample  $t$  procedures and analysis of variance) are quite **robust**. That is, the results of inference are not very sensitive to moderate lack of Normality, especially when the samples are reasonably large. Some practical guidelines for taking advantage of the robustness of these methods appear in Chapter 7.

robustness

What can we do if plots suggest that the population distribution is clearly not Normal, especially when we have only a few observations? This is not a simple question. Here are the basic options:

- outlier** 1. If lack of Normality is due to **outliers**, it may be legitimate to remove the outliers. An outlier is an observation that may not come from the same population as the other observations. Equipment failure that produced a bad measurement, for example, entitles you to remove the outlier and analyze the remaining data. If the outlier appears to be “real data,” you can base inference on statistics that are more resistant than  $\bar{x}$  and  $s$ . Options 4 and 5 allow this.
- 2. Sometimes we can **transform** our data so that their distribution is more nearly Normal. Transformations such as the logarithm that pull in the long tail of right-skewed distributions are particularly helpful. Example 7.10 (page 421) illustrates use of the logarithm. A detailed discussion of transformations appears in the extra material entitled *Transforming Relationships* available on the course Web site.
- other standard distributions** 3. In some settings, **other standard distributions** replace the Normal distributions as models for the overall pattern in the population. We mentioned in Chapter 5 (page 330) that the Weibull distributions are common models for the lifetimes in service of equipment in statistical studies of reliability. There are inference procedures for the parameters of these distributions that replace the  $t$  procedures when we use specific non-Normal models.
- bootstrap methods**  
**permutation tests** 4. Modern **bootstrap methods** and **permutation tests** do not require Normality or any other specific form of sampling distribution. Moreover, you can base inference on resistant statistics such as the trimmed mean. We recommend these methods unless the sample is so small that it may not represent the population well. Chapter 16 gives a full discussion.
- nonparametric methods** 5. Finally, there are other **nonparametric methods** that do not require any specific form for the distribution of the population. Unlike bootstrap and permutation methods, common nonparametric methods do not make use of the actual values of the observations. The *sign test* (page 423) works with *counts* of observations. This chapter presents **rank tests** based on the *rank* (place in order) of each observation in the set of all the data.
- rank tests**

This chapter concerns rank tests that are designed to replace the  $t$  tests and one-way analysis of variance when the Normality conditions for those tests are not met. Figure 15.1 presents an outline of the standard tests (based on Normal distributions) and the rank tests that compete with them. All these tests require that the population or populations have **continuous distributions**. That is, each distribution must be described by a density curve that allows observations to take any value in some interval of outcomes. The Normal curves are one shape of density curve. Rank tests allow curves of any shape.

The rank tests we will study concern the *center* of a population or populations. When a population has at least roughly a Normal distribution, we describe its center by the mean. The “Normal tests” in Figure 15.1 test hypotheses about population means. When distributions are strongly skewed, we often prefer the median to the mean as a measure of center. In simplest form, the hypotheses for rank tests just replace mean by median.

 **LOOK BACK**  
continuous distribution  
p. 253

**FIGURE 15.1** Comparison of tests based on Normal distributions with nonparametric tests for similar settings.

Setting	Normal test	Rank test
One sample	One-sample $t$ test Section 7.1	Wilcoxon signed rank test Section 15.2
Matched pairs	Apply one-sample test to differences within pairs	
Two independent samples	Two-sample $t$ test Section 7.2	Wilcoxon rank sum test Section 15.1
Several independent samples	One-way ANOVA $F$ test Chapter 12	Kruskal-Wallis test Section 15.3

We devote a section of this chapter to each of the rank procedures. Section 15.1, which discusses the most common of these tests, also contains general information about rank tests. The kind of assumptions required, the nature of the hypotheses tested, the big idea of using ranks, and the contrast between exact distributions for use with small samples and approximations for use with larger samples are common to all rank tests. Sections 15.2 and 15.3 more briefly describe other rank tests.

## 15.1 The Wilcoxon Rank Sum Test

← **LOOK BACK**  
two sample problems  
p. 432

Two-sample problems (see Section 7.2) are among the most common in statistics. The most useful nonparametric significance test compares two distributions. Here is an example of this setting.

### EXAMPLE

**15.1 Weeds and corn yield.** Does the presence of small numbers of weeds reduce the yield of corn? Lamb’s-quarter is a common weed in corn fields. A researcher planted corn at the same rate in eight small plots of ground, then weeded the corn rows by hand to allow no weeds in four randomly selected plots and exactly three lamb’s-quarter plants per meter of row in the other four plots. Here are the yields of corn (bushels per acre) in each of the plots:<sup>1</sup>

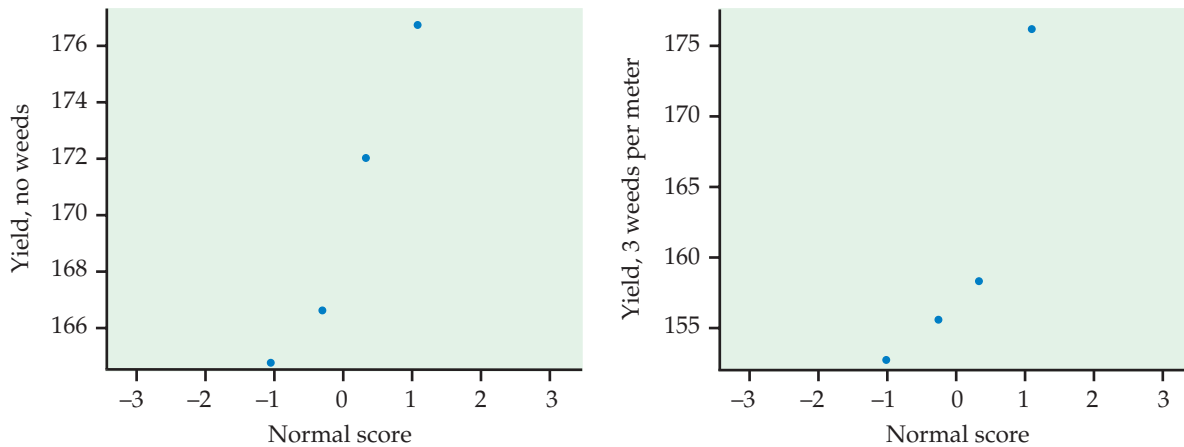
Weeds per meter	Yield (bu/acre)			
0	166.7	172.2	165.0	176.9
3	158.6	176.4	153.1	156.0

Normal quantile plots (Figure 15.2) suggest that the population distribution may be right-skewed. The samples are too small to assess Normality adequately or to rely on the robustness of the two-sample  $t$  test. We prefer to use a test that does not require Normality.

### The rank transformation

We first rank all eight observations together. To do this, arrange them in order from smallest to largest:

153.1 156.0 158.6 **165.0** **166.7** **172.2** 176.4 **176.9**



**FIGURE 15.2** Normal quantile plots of corn yields from plots with no weeds (left) and with three weeds per meter of row (right), for Example 15.1.

The boldface entries in the list are the yields with no weeds present. We see that four of the five highest yields come from that group, suggesting that yields are higher with no weeds. The idea of rank tests is to look just at position in this ordered list. To do this, replace each observation by its order, from 1 (smallest) to 8 (largest). These numbers are the *ranks*:

Yield	153.1	156.0	158.6	<b>165.0</b>	<b>166.7</b>	<b>172.2</b>	176.4	<b>176.9</b>
Rank	1	2	3	<b>4</b>	<b>5</b>	<b>6</b>	7	<b>8</b>

**RANKS**

To rank observations, first arrange them in order from smallest to largest. The **rank** of each observation is its position in this ordered list, starting with rank 1 for the smallest observation.

Moving from the original observations to their ranks is a transformation of the data, like moving from the observations to their logarithms. The rank transformation retains only the ordering of the observations and makes no other use of their numerical values. Working with ranks allows us to dispense with specific assumptions about the shape of the distribution, such as Normality.

**USE YOUR KNOWLEDGE**

**15.1 Numbers of rooms in top spas.** A report of a readers' poll in *Condé Nast Traveler* magazine ranked 36 top resort spas.<sup>2</sup> Let Group A be the top-ranked 18 spas, and let Group B be the next 18 rated spas in the list. A simple random sample of size 5 was taken from each group, and the number of rooms in each selected spa was recorded. Here are the data:

Group A	552	448	68	243	30
Group B	329	780	560	540	240





Rank all of the observations together and make a list of the ranks for Group A and Group B.

**15.2 The effect of Spa Bellagio on the result.** Refer to the previous exercise. Spa Bellagio in Las Vegas is one of the spas in Group B. Suppose this spa had been the second spa selected in the random sample for Group B. Replace the observation 780 in Group B by 4003, the number of rooms in Spa Bellagio. Use the modified data to make a list of the ranks for Groups A and B combined. What changes?

### The Wilcoxon rank sum test

If the presence of weeds reduces corn yields, we expect the ranks of the yields from plots with weeds to be smaller as a group than the ranks from plots without weeds. We might compare the *sums* of the ranks from the two treatments:

Treatment	Sum of ranks
No weeds	23
Weeds	13

These sums measure how much the ranks of the weed-free plots as a group exceed those of the weedy plots. In fact, the sum of the ranks from 1 to 8 is always equal to 36, so it is enough to report the sum for one of the two groups. If the sum of the ranks for the weed-free group is 23, the ranks for the other group must add to 13 because  $23 + 13 = 36$ . If the weeds have no effect, we would expect the sum of the ranks in each group to be 18 (half of 36). Here are the facts we need in a more general form that takes account of the fact that our two samples need not be the same size.

#### THE WILCOXON RANK SUM TEST

Draw an SRS of size  $n_1$  from one population and draw an independent SRS of size  $n_2$  from a second population. There are  $N$  observations in all, where  $N = n_1 + n_2$ . Rank all  $N$  observations. The sum  $W$  of the ranks for the first sample is the **Wilcoxon rank sum statistic**. If the two populations have the same continuous distribution, then  $W$  has mean

$$\mu_W = \frac{n_1(N + 1)}{2}$$

and standard deviation

$$\sigma_W = \sqrt{\frac{n_1 n_2 (N + 1)}{12}}$$

The **Wilcoxon rank sum test** rejects the hypothesis that the two populations have identical distributions when the rank sum  $W$  is far from its mean.\*

\*This test was invented by Frank Wilcoxon (1892–1965) in 1945. Wilcoxon was a chemist who encountered statistical problems in his work at the research laboratories of American Cyanamid Company.

In the corn yield study of Example 15.1, we want to test

$H_0$ : no difference in distribution of yields

against the one-sided alternative

$H_a$ : yields are systematically higher in weed-free plots

Our test statistic is the rank sum  $W = 23$  for the weed-free plots.



### USE YOUR KNOWLEDGE

**15.3 Hypotheses and test statistic for top spas.** Refer to Exercise 15.1. State appropriate null and alternative hypotheses for this setting and calculate the value of  $W$ , the test statistic.

**15.4 Effect of Spa Bellagio on the test statistic.** Refer to Exercise 15.2. Using the altered data, state appropriate null and alternative hypotheses and calculate the value of  $W$ , the test statistic.

### EXAMPLE

**15.2 Perform the significance test.** In Example 15.1,  $n_1 = 4$ ,  $n_2 = 4$ , and there are  $N = 8$  observations in all. The sum of ranks for the weed-free plots has mean

$$\begin{aligned}\mu_W &= \frac{n_1(N+1)}{2} \\ &= \frac{(4)(9)}{2} = 18\end{aligned}$$

and standard deviation

$$\begin{aligned}\sigma_W &= \sqrt{\frac{n_1 n_2 (N+1)}{12}} \\ &= \sqrt{\frac{(4)(4)(9)}{12}} = \sqrt{12} = 3.464\end{aligned}$$

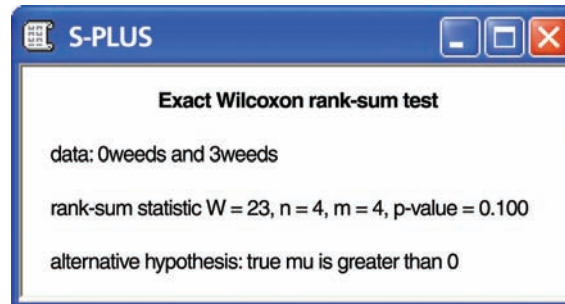
Although the observed rank sum  $W = 23$  is higher than the mean, it is only about 1.4 standard deviations higher. We now suspect that the data do not give strong evidence that yields are higher in the population of weed-free corn.

The  $P$ -value for our one-sided alternative is  $P(W \geq 23)$ , the probability that  $W$  is at least as large as the value for our data when  $H_0$  is true.

To calculate the  $P$ -value  $P(W \geq 23)$ , we need to know the sampling distribution of the rank sum  $W$  when the null hypothesis is true. This distribution depends on the two sample sizes  $n_1$  and  $n_2$ . Tables are therefore a bit unwieldy, though you can find them in handbooks of statistical tables. Most statistical software will give you  $P$ -values, as well as carry out the ranking and calculate  $W$ . However, some software gives only approximate  $P$ -values. You must learn what your software offers.

**EXAMPLE**

**15.3 Software output.** Figure 15.3 shows the output from software that calculates the exact sampling distribution of  $W$ . We see that the sum of the ranks in the weed-free group is  $W = 23$ , with  $P$ -value  $P = 0.100$  against the one-sided alternative that weed-free plots have higher yields. There is some evidence that weeds reduce yield, considering that we have data from only four plots for each treatment. The evidence does not, however, reach the levels usually considered convincing.



**FIGURE 15.3** Output from the S-PLUS statistical software for the data in Example 15.1. The program uses the exact distribution for  $W$  when the samples are small and there are no tied observations.

← **LOOK BACK**  
two-sample  $t$  test  
p. 436

It is worth noting that the two-sample  $t$  test gives essentially the same result as the Wilcoxon test in Example 15.3 ( $t = 1.554$ ,  $P = 0.0937$ ). A permutation test (Chapter 16) for the sample means gives  $P = 0.084$ . It is in fact somewhat unusual to find a strong disagreement among the conclusions reached by these tests.

### The Normal approximation

The rank sum statistic  $W$  becomes approximately Normal as the two sample sizes increase. We can then form yet another  $z$  statistic by standardizing  $W$ :

$$\begin{aligned} z &= \frac{W - \mu_W}{\sigma_W} \\ &= \frac{W - n_1(N + 1)/2}{\sqrt{n_1 n_2 (N + 1)/12}} \end{aligned}$$

← **LOOK BACK**  
continuity correction,  
p. 327

Use standard Normal probability calculations to find  $P$ -values for this statistic. Because  $W$  takes only whole-number values, the **continuity correction** improves the accuracy of the approximation.

**EXAMPLE**

**15.4 The continuity correction.** The standardized rank sum statistic  $W$  in our corn yield example is

$$z = \frac{W - \mu_W}{\sigma_W} = \frac{23 - 18}{3.464} = 1.44$$

We expect  $W$  to be larger when the alternative hypothesis is true, so the approximate  $P$ -value is

$$P(Z \geq 1.44) = 0.0749$$

The continuity correction acts as if the whole number 23 occupies the entire interval from 22.5 to 23.5. We calculate the  $P$ -value  $P(W \geq 23)$  as  $P(W \geq 22.5)$  because the value 23 is included in the range whose probability we want. Here is the calculation:

$$\begin{aligned} P(W \geq 22.5) &= P\left(\frac{W - \mu_W}{\sigma_W} \geq \frac{22.5 - 18}{3.464}\right) \\ &= P(Z \geq 1.30) \\ &= 0.0968 \end{aligned}$$

The continuity correction gives a result closer to the exact value  $P = 0.100$ .



## USE YOUR KNOWLEDGE

**15.5 The  $P$ -value for top spas.** Refer to Exercises 15.1 and 15.3. Find  $\mu_W$ ,  $\sigma_W$ , and the standardized rank sum statistic. Then give an approximate  $P$ -value using the Normal approximation. What do you conclude?

**15.6 The effect of Spa Bellagio on the  $P$ -value.** Refer to Exercises 15.2 and 15.4. Answer the questions for Exercise 15.5 using the altered data.

We recommend always using either the exact distribution (from software or tables) or the continuity correction for the rank sum statistic  $W$ . The exact distribution is safer for small samples. As Example 15.4 illustrates, however, the Normal approximation with the continuity correction is often adequate.

## EXAMPLE

**15.5 Software output.** Figure 15.4 shows the output for our data from two more statistical programs. Minitab offers only the Normal approximation, and it refers to the **Mann-Whitney test**. This is an alternative form of the Wilcoxon rank sum test. SAS carries out both the exact and the approximate tests. SAS calls the rank sum  $S$  rather than  $W$  and gives the mean 18 and standard deviation 3.464 as well as the  $z$  statistic 1.299 (using the continuity correction). SAS gives the approximate two-sided  $P$ -value as 0.1939, so the one-sided result is half this,  $P = 0.0970$ . This agrees with Minitab and (up to a small roundoff error) with our result in Example 15.4. This approximate  $P$ -value is close to the exact result  $P = 0.100$ , given by SAS and in Figure 15.4.

Mann-Whitney test

## What hypotheses does Wilcoxon test?

Our null hypothesis is that weeds do not affect yield. Our alternative hypothesis is that yields are lower when weeds are present. If we are willing to assume that yields are Normally distributed, or if we have reasonably large samples, we use the two-sample  $t$  test for means. Our hypotheses then become

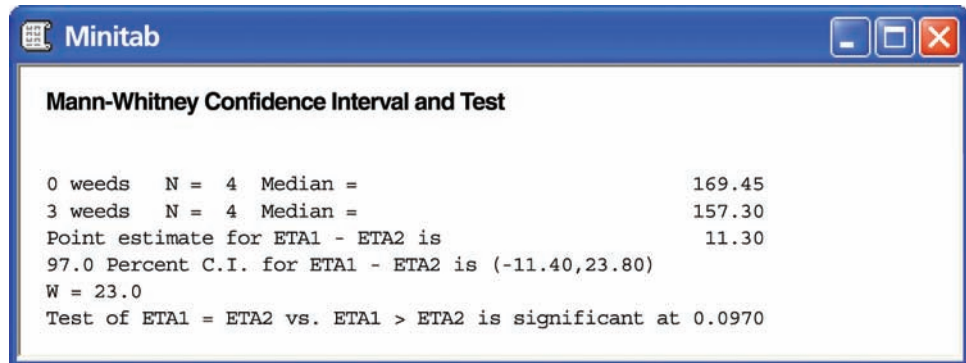
$$H_0: \mu_1 = \mu_2$$

$$H_a: \mu_1 > \mu_2$$

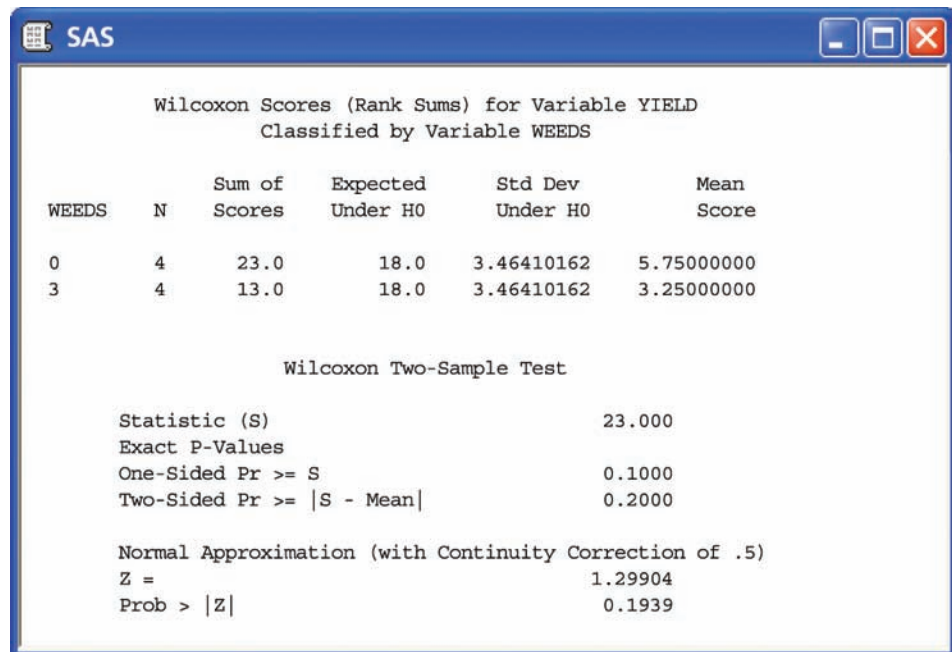
When the distributions may not be Normal, we might restate the hypotheses in terms of population medians rather than means:

$$H_0: \text{median}_1 = \text{median}_2$$

$$H_a: \text{median}_1 > \text{median}_2$$



(a)



(b)

**FIGURE 15.4** Output from the Minitab and SAS statistical software for the data in Example 15.1. (a) Minitab uses the Normal approximation for the distribution of  $W$ . (b) SAS gives both the exact and approximate values.



*The Wilcoxon rank sum test does test hypotheses about population medians, but only if an additional assumption is met: both populations must have distributions of the same shape. That is, the density curve for corn yields with three weeds per meter looks exactly like that for no weeds except that it may slide to a different location on the scale of yields. The Minitab output in Figure 15.4(a) states the hypotheses in terms of population medians (which it calls “ETA”) and also gives a confidence interval for the difference between the two population medians.*

The same-shape assumption is too strict to be reasonable in practice. Recall that our preferred version of the two-sample  $t$  test does not require that the two populations have the same standard deviation—that is, it does not make a same-shape assumption. Fortunately, the Wilcoxon test also applies in a much

more general and more useful setting. It tests hypotheses that we can state in words as

$H_0$ : The two distributions are the same.

$H_a$ : One distribution has values that are systematically larger.

Here is a more exact statement of the “systematically larger” alternative hypothesis. Take  $X_1$  to be corn yield with no weeds and  $X_2$  to be corn yield with three weeds per meter. These yields are random variables. That is, every time we plant a plot with no weeds, the yield is a value of the variable  $X_1$ . The probability that the yield is more than 160 bushels per acre when no weeds are present is  $P(X_1 > 160)$ . If weed-free yields are “systematically larger” than those with weeds, yields higher than 160 should be more likely with no weeds. That is, we should have

$$P(X_1 > 160) > P(X_2 > 160)$$

The alternative hypothesis says that this inequality holds not just for 160 but for *any* yield we care to specify. No weeds always puts more probability “to the right” of whatever yield we are interested in.<sup>3</sup>

This exact statement of the hypotheses we are testing is a bit awkward. The hypotheses really are “nonparametric” because they do not involve any specific parameter such as the mean or median. If the two distributions do have the same shape, the general hypotheses reduce to comparing medians. Many texts and computer outputs state the hypotheses in terms of medians, sometimes ignoring the same-shape requirement. We recommend that you express the hypotheses in words rather than symbols. “Yields are systematically higher in weed-free plots” is easy to understand and is a good statement of the effect that the Wilcoxon test looks for.

## Ties

The exact distribution for the Wilcoxon rank sum is obtained assuming that all observations in both samples take different values. This allows us to rank them all. In practice, however, we often find observations tied at the same value. What shall we do? The usual practice is to *assign all tied values the average of the ranks they occupy*. Here is an example with six observations:

average ranks

Observation	153	155	158	158	161	164
Rank	1	2	3.5	3.5	5	6

The tied observations occupy the third and fourth places in the ordered list, so they share rank 3.5.

The exact distribution for the Wilcoxon rank sum  $W$  changes if the data contain ties. Moreover, the standard deviation  $\sigma_W$  must be adjusted if ties are present. The Normal approximation can be used after the standard deviation is adjusted. Statistical software will detect ties, make the necessary adjustment, and switch to the Normal approximation. In practice, software is required if you want to use rank tests when the data contain tied values.

It is sometimes useful to use rank tests on data that have very many ties because the scale of measurement has only a few values. Here is an example.

**EXAMPLE**

**15.6 Job satisfaction.** Self-employed people generally have more control over their work than those who work for others. Does this difference translate into greater job satisfaction? A Pew Research Center survey compared the job satisfaction rating of workers who were self-employed with those who are not.<sup>4</sup> Here are the responses:

	Count				Total
	Completely Satisfied	Mostly Satisfied	Mostly Dissatisfied	Completely Dissatisfied	
Self-employed	99	142	8	5	254
Not self-employed	250	542	73	20	885

**USE YOUR KNOWLEDGE**

**15.7 Analyze as a two-way table.** Use the data in Example 15.6.

- Compute the percents of the different responses for the self-employed workers. Do the same for those who are not self-employed. Display the percents graphically and summarize the differences in the two distributions.
- Perform the chi-square test for the counts in the two-way table. Report the test statistic, the degrees of freedom, and the  $P$ -value. Give a brief summary of what you can conclude from this significance test.

How do we approach the analysis of these data using the Wilcoxon test? We start with the hypotheses. We have two distributions of job satisfaction, one for those who are self-employed and one for those who are not. The null hypothesis states that the two distributions are the same. The alternative hypothesis uses the fact that the responses are ordered from the most satisfied to the least satisfied. It states that one of the employment groups is more satisfied than the other.

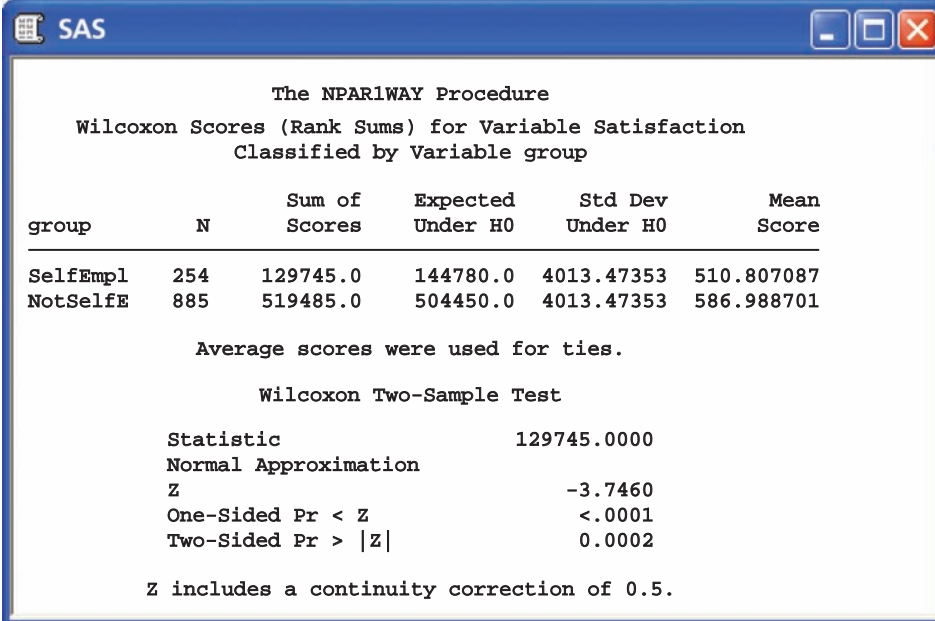
$H_0$ : Self-employed workers and those who are not self employed have the same job satisfaction.

$H_a$ : One of the two groups of workers has greater job satisfaction than the other.

The alternative hypothesis is two-sided. Because the responses can take only four values, there are very many ties. All 25 workers who are completely dissatisfied are tied. Similarly, all workers in each of the four columns of the table corresponding to the different responses are tied. The graphical display that you prepared in Exercise 15.7 suggests that self-employed workers have greater job satisfaction. Is this difference statistically significant?

**EXAMPLE**

**15.7 Software output.** Look at Figure 15.5, which gives software output for the Wilcoxon test. The rank sum for the self-employed workers (using average ranks for ties) is  $W = 129,745$ . The standardized value for this statistic is  $z = -3.75$  and the two-sided  $P$ -value is  $P = 0.0002$ . There is very strong evidence of a difference. Self-employed workers have greater job satisfaction than workers who are not self-employed.



The NPAR1WAY Procedure  
Wilcoxon Scores (Rank Sums) for Variable Satisfaction  
Classified by Variable group

group	N	Sum of Scores	Expected Under H0	Std Dev Under H0	Mean Score
SelfEmpl	254	129745.0	144780.0	4013.47353	510.807087
NotSelfE	885	519485.0	504450.0	4013.47353	586.988701

Average scores were used for ties.

Wilcoxon Two-Sample Test

Statistic	129745.0000
Normal Approximation	
Z	-3.7460
One-Sided Pr < Z	<.0001
Two-Sided Pr >  Z	0.0002

Z includes a continuity correction of 0.5.

**FIGURE 15.5** Output from SAS for the job satisfaction survey of Example 15.6. The approximate two-sided  $P$ -value is 0.0002.

With more than 200 observations in each group and no outliers, we might use the two-sample  $t$  test (even though responses take only four values). To perform the  $t$  test, we recode the four responses numerically, using the values 1, 2, 3, and 4 for the responses “completely satisfied,” “mostly satisfied,” “mostly dissatisfied,” and “completely dissatisfied.” The results are  $t = 3.54$  with  $P = 0.0004$ . The  $P$ -value for two-sample  $t$  test is essentially the same as that for the Wilcoxon test. There is, however, another reason to prefer the rank test in this example. The  $t$  statistic treats the response values 1 through 5 as meaningful numbers. In particular, the possible responses are treated as though they are equally spaced. The difference between “completely satisfied” and “mostly satisfied” is the same as the difference between “mostly satisfied” and “mostly dissatisfied.” This may not make sense. The rank test, on the other hand, uses only the order of the responses, not their actual values. The responses are arranged in order from most satisfied to least satisfied, so the rank test makes sense. *Some statisticians avoid using  $t$  procedures when there is not a fully meaningful scale of measurement.*



### Rank, $t$ , and permutation tests

The two-sample  $t$  procedures are the most common method for comparing the centers of two populations based on random samples from each. The Wilcoxon

rank sum test is a competing procedure that does not start from the condition that the populations have Normal distributions. Permutation tests (Chapter 16) also avoid the need for Normality. Tests based on Normality, rank tests, and permutation tests apply in many other settings as well. How do these three approaches compare in general?

First consider rank tests versus traditional tests based on Normal distributions. Both are available in almost all statistical software.

- Moving from the actual data values to their ranks allows us to find an exact sampling distribution for rank statistics such as the Wilcoxon rank sum  $W$  when the null hypothesis is true. (Most software will do this only if there are no ties and if the samples are quite small.) When our samples are small, are truly random samples from the populations, and show non-Normal distributions of the same shape, the Wilcoxon test is more reliable than the two-sample  $t$  test. In practice, the robustness of  $t$  procedures implies that we rarely encounter data that require nonparametric procedures to obtain reasonably accurate  $P$ -values. The  $t$  and  $W$  tests give very similar results in our examples. Nonetheless, many statisticians would not use a  $t$  test in Example 15.6 because the response variable gives only the order of the responses.
- Normal tests compare means and are accompanied by simple confidence intervals for means or differences between means. When we use rank tests to compare medians, we can also give confidence intervals for medians. However, the usefulness of rank tests is clearest in settings when they do not simply compare medians—see the discussion “What hypotheses does Wilcoxon test?” Rank methods emphasize tests, not confidence intervals.
- Inference based on ranks is largely restricted to simple settings. Normal inference extends to methods for use with complex experimental designs and multiple regression, but nonparametric tests do not. We stress Normal inference in part because it leads to more advanced statistics.

If you have already read Chapter 16 and use software that makes permutation tests available to you, you will also want to compare rank tests with resampling methods.

- Both rank and permutation tests are nonparametric. That is, they require no assumptions about the shape of the population distribution. A two-sample permutation test has the same null hypothesis as the Wilcoxon rank sum test: that the two population distributions are identical. Calculation of the sampling distribution under the null hypothesis is similar for both tests but is simpler for rank tests because it depends only on the sizes of the samples. As a result, software often gives exact  $P$ -values for rank tests but not for permutation tests.
- Permutation tests have the advantage of flexibility. They allow wide choice of the statistic used to compare two samples, an advantage over both the  $t$  and Wilcoxon tests. In fact, we could apply the permutation test method to sample means (imitating  $t$ ) or to rank sums (imitating Wilcoxon), as well as to other statistics such as the trimmed mean. Permutation tests are not available in some settings, such as testing hypotheses about a single population, though bootstrap confidence intervals do allow resampling tests in these settings. Permutation tests are available for multiple regression and some other quite elaborate settings.



**LOOK BACK**  
trimmed mean  
p. 199

- An important advantage of resampling methods over both Normal and rank procedures is that we can get bootstrap confidence intervals for the parameter corresponding to whatever statistic we choose for the permutation test. If the samples are very small, however, bootstrap confidence intervals may be unreliable because the samples don't represent the population well enough to provide a good basis for bootstrapping.



In general, both Normal distribution methods and resampling methods are more useful than rank tests. *If you are familiar with resampling, we recommend rank tests only for very small samples, and even then only if your software gives exact P-values for rank tests but not for permutation tests.*

## SECTION 15.1 Summary

**Nonparametric tests** do not require any specific form for the distribution of the population from which our samples come.

**Rank tests** are nonparametric tests based on the **ranks** of observations, their positions in the list ordered from smallest (rank 1) to largest. Tied observations receive the average of their ranks.


The **Wilcoxon rank sum test** compares two distributions to assess whether one has systematically larger values than the other. The Wilcoxon test is based on the **Wilcoxon rank sum statistic  $W$** , which is the sum of the ranks of one of the samples. The Wilcoxon test can replace the **two-sample  $t$  test**.

**P-values** for the Wilcoxon test are based on the sampling distribution of the rank sum statistic  $W$  when the null hypothesis (no difference in distributions) is true. You can find  $P$ -values from special tables, software, or a Normal approximation (with continuity correction).

## SECTION 15.1 Exercises

For Exercises 15.1 and 15.2, see pages 15-4 to 15-5; for Exercises 15.3 and 15.4, see page 15-6; for Exercises 15.5 and 15.6, see page 15-8; and for Exercise 15.7, see page 15-11.


Statistical software is very helpful in doing these exercises. If you do not have access to software, base your work on the Normal approximation with continuity correction (page 15-7).

**15.8 Do women talk more?** Conventional wisdom suggests that women are more talkative than men. One study designed to examine this stereotype collected data on the speech of 10 men and 10 women in the United States.<sup>5</sup> The variable recorded is the number of words per day. Here are the data:  TALK10

Men				Women			
23871	5180	9951	12460	10592	24608	13739	22376
17155	10344	9811	12387	9351	7694	16812	21066
29920	21791			32291	12320		

(a) Summarize the data for the two groups using numerical and graphical methods. Describe the two distributions.

(b) Compare the words per day spoken by the men with the words per day spoken by the women using the Wilcoxon rank sum test. Summarize your results and conclusion in a short paragraph.


**15.9 More data for women and men talking.** The data in the previous exercise were a sample of the data collected in a larger study of 42 men and 37 women. Use the larger data set to answer the questions in the previous exercise. Discuss the advisability of using the Wilcoxon test versus the  $t$  test for this exercise and for the previous one.  TALK

**15.10 Weeds and corn yield.** The corn yield study of Example 15.1 also examined yields in four plots having nine lamb's-quarter plants per meter of row. The yields (bushels per acre) in these plots were

162.8   142.4   162.7   162.4

There is a clear outlier, but rechecking the results found that this is the correct yield for this plot. The outlier makes us hesitant to use  $t$  procedures because  $\bar{x}$  and  $s$  are not resistant.

- (a) Is there evidence that 9 weeds per meter reduces corn yields when compared with weed-free corn? Use the Wilcoxon rank sum test with the preceding data and some of the data from Example 15.1 to answer this question.
- (b) Compare the results from part (a) with those from the two-sample  $t$  test for these data.
- (c) Now remove the low outlier 142.4 from the data for nine weeds per meter. Repeat both the Wilcoxon and  $t$  analyses. By how much did the outlier reduce the mean yield in its group? By how much did it increase the standard deviation? Did it have a practically important impact on your conclusions?

**15.11 Storytelling and the use of language.** A study of early childhood education asked kindergarten students to retell two fairy tales that had been read to them earlier in the week. The 10 children in the study included 5 high-progress readers and 5 low-progress readers. Each child told two stories. Story 1 had been read to them; Story 2 had been read and also illustrated with pictures. An expert listened to a recording of each child and assigned a score for certain uses of language. Here are the data.<sup>6</sup>  STORYTELLING

Child	Progress	Story 1 score	Story 2 score	Child	Progress	Story 1 score	Story 2 score
1	high	0.55	0.80	6	low	0.40	0.77
2	high	0.57	0.82	7	low	0.72	0.49
3	high	0.72	0.54	8	low	0.00	0.66
4	high	0.70	0.79	9	low	0.36	0.28
5	high	0.84	0.89	10	low	0.55	0.38

Is there evidence that the scores of high-progress readers are higher than those of low-progress readers when they retell a story they have heard without pictures (Story 1)?

- (a) Make Normal quantile plots for the 5 responses in each group. Are any major deviations from Normality apparent?
- (b) Carry out a two-sample  $t$  test. State hypotheses and give the two sample means, the  $t$  statistic and its  $P$ -value, and your conclusion.
- (c) Carry out the Wilcoxon rank sum test. State hypotheses and give the rank sum  $W$  for high-progress readers, its  $P$ -value, and your conclusion. Do the  $t$  and Wilcoxon tests lead you to different conclusions?


**15.12 Repeat the analysis for Story 2.** Repeat the analysis of Exercise 15.11 for the scores when children

retell a story they have heard and seen illustrated with pictures (Story 2).  STORYTELLING

**15.13 Do the calculations by hand.** Use the data in Exercise 15.11 for children telling Story 2 to carry out by hand the steps in the Wilcoxon rank sum test.

 STORYTELLING

- (a) Arrange the 10 observations in order and assign ranks. There are no ties.
- (b) Find the rank sum  $W$  for the 5 high-progress readers. What are the mean and standard deviation of  $W$  under the null hypothesis that low-progress and high-progress readers do not differ?
- (c) Standardize  $W$  to obtain a  $z$  statistic. Do a Normal probability calculation with the continuity correction to obtain a one-sided  $P$ -value.
- (d) The data for Story 1 contain tied observations. What ranks would you assign to the 10 scores for Story 1?


**15.14 Learning math through subliminal messages.** A “subliminal” message is below our threshold of awareness but may nonetheless influence us. Can subliminal messages help students learn math? A group of students who had failed the mathematics part of the City University of New York Skills Assessment Test agreed to participate in a study to find out. All received a daily subliminal message, flashed on a screen too rapidly to be consciously read. The treatment group of 10 students was exposed to “Each day I am getting better in math.” The control group of 8 students was exposed to a neutral message, “People are walking on the street.” All students participated in a summer program designed to raise their math skills, and all took the assessment test again at the end of the program. Here are data on the subjects’ scores before and after the program.<sup>7</sup>  SUBLIMINALMATH

Treatment Group		Control Group	
Pretest	Posttest	Pretest	Posttest
18	24	18	29
18	25	24	29
21	33	20	24
18	29	18	26
18	33	24	38
20	36	22	27
23	34	15	22
23	36	19	31
21	34		
17	27		

- (a) The study design was a randomized comparative experiment. Outline this design.
- (b) Compare the gain in scores in the two groups, using a graph and numerical descriptions. Does it appear that the

treatment group's scores rose more than the scores for the control group?


(c) Apply the Wilcoxon rank sum test to the posttest versus pretest differences. Note that there are some ties. What do you conclude?

**15.15 Effects of logging in Borneo.** "Conservationists have despaired over destruction of tropical rainforest by logging, clearing, and burning." These words begin a report on a statistical study of the effects of logging in Borneo.<sup>8</sup> Here are data on the number of tree species in 12 unlogged forest plots and 9 similar plots logged eight years earlier:  BORNEO

Unlogged		Logged	
22	18	17	4
22	20	18	14
15	21	18	15
13	13	15	10
19	13	12	
19	15		


(a) Make a back-to-back stemplot of the data. Does there appear to be a difference in species counts for logged and unlogged plots?

(b) Does logging significantly reduce the number of species in a plot after eight years? State hypotheses, do a Wilcoxon test, and state your conclusion.

**15.16 Improved methods for teaching reading.** Do new "directed reading activities" improve the reading ability of elementary school students, as measured by their Degree of Reading Power (DRP) score? A study assigns students at random to either the new method (treatment group, 21 students) or traditional teaching methods (control group, 23 students). Here are the DRP scores at the end of the study:<sup>9</sup>  READINGDRP

Treatment group				Control group			
24	61	59	46	42	33	46	37
43	44	52	43	43	41	10	42
58	67	62	57	55	19	17	55
71	49	54		26	54	60	28
43	53	57		62	20	53	48
49	56	33		37	85	42	

For these data the two-sample  $t$  test (Example 7.14, page 436) gives  $P = 0.013$  and a permutation test based on the difference of means (Example 16.12, page 16-43) gives  $P = 0.015$ . Both of these tests are based on the difference of sample means. Does the Wilcoxon test, based on rank sums rather than means, give a similar  $P$ -value?

**15.17 Attitudes toward secondhand stores.** To study customers' attitudes toward secondhand stores, researchers interviewed samples of shoppers at two secondhand stores of the same chain in two cities. Here are data on the incomes of shoppers at the two stores, presented as a two-way table of counts:<sup>10</sup>  SECONDHAND

Income	City 1	City 2
Under \$10,000	70	62
\$10,000 to \$19,999	52	63
\$20,000 to \$24,999	69	50
\$25,000 to \$34,999	22	19
\$35,000 or more	28	24

(a) Is there a relationship between city and income? Use the chi-square test to answer this question.

(b) The chi-square test ignores the ordering of the income categories. Is there good evidence that shoppers in one city have systematically higher incomes than in the other?

## 15.2 The Wilcoxon Signed Rank Test

We use the one-sample  $t$  procedures for inference about the mean of one population or for inference about the mean difference in a matched pairs setting. The matched pairs setting is more important because good studies are generally comparative. We will now meet a rank test for this setting.

### EXAMPLE

**15.8 Storytelling and reading.** A study of early childhood education asked kindergarten students to retell two fairy tales that had been read to them earlier in the week. Each child told two stories. The first had been read to them, and the second had been read but also illustrated with pictures. An expert



listened to a recording of the children and assigned a score for certain uses of language. Here are the data for five “low-progress” readers in a pilot study:<sup>11</sup>

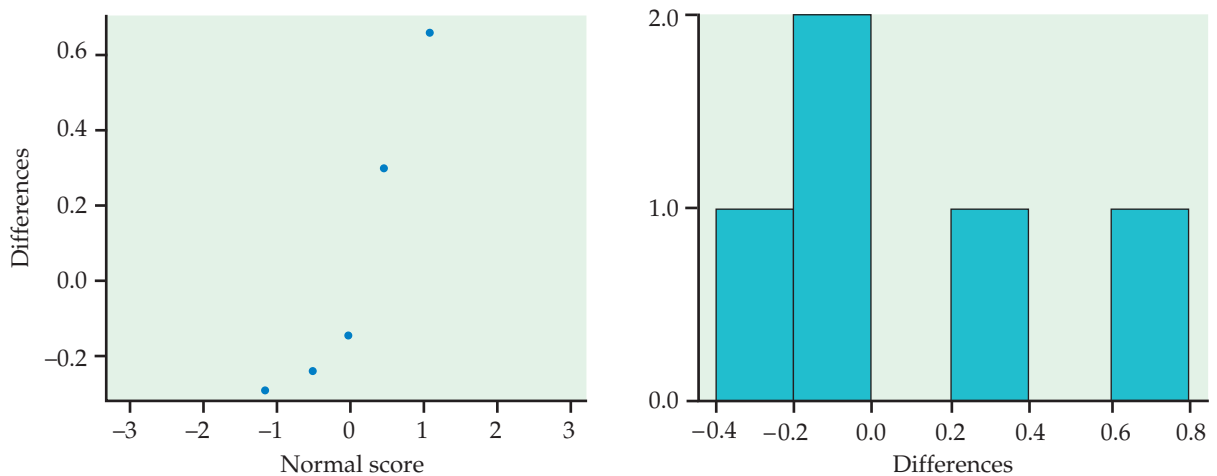
Child	1	2	3	4	5
Story 2	0.77	0.49	0.66	0.28	0.38
Story 1	0.40	0.72	0.00	0.36	0.55
Difference	0.37	-0.23	0.66	-0.08	-0.17

We wonder if illustrations improve how the children retell a story. We would like to test the hypotheses

$H_0$ : Scores have the same distribution for both stories.

$H_a$ : Scores are systematically higher for Story 2.

Because this is a matched pairs design, we base our inference on the differences. The matched pairs  $t$  test gives  $t = 0.635$  with one-sided  $P$ -value  $P = 0.280$ . Displays of the data (Figure 15.6) suggest some lack of Normality. We would therefore like to use a rank test.



**FIGURE 15.6** Normal quantile plot and histogram for the five differences in Example 15.8.

Positive differences in Example 15.8 indicate that the child performed better telling Story 2. If scores are generally higher with illustrations, the positive differences should be farther from zero in the positive direction than the negative differences are in the negative direction. We therefore compare the **absolute values** of the differences, that is, their magnitudes without a sign. Here they are, with boldface indicating the positive values:

absolute value

**0.37** 0.23 **0.66** 0.08 0.17

Arrange these in increasing order and assign ranks, keeping track of which values were originally positive. Tied values receive the average of their ranks. If there are cases with zero differences, discard them before ranking.

Absolute value	0.08	0.17	0.23	<b>0.37</b>	<b>0.66</b>
Rank	1	2	3	<b>4</b>	<b>5</b>

The test statistic is the sum of the ranks of the positive differences. (We could equally well use the sum of the ranks of the negative differences.) This is the *Wilcoxon signed rank statistic*. Its value here is  $W^+ = 9$ .

### THE WILCOXON SIGNED RANK TEST FOR MATCHED PAIRS

Draw an SRS of size  $n$  from a population for a matched pairs study and take the differences in responses within pairs. Rank the absolute values of these differences. The sum  $W^+$  of the ranks for the positive differences is the **Wilcoxon signed rank statistic**. If the distribution of the responses is not affected by the different treatments within pairs, then  $W^+$  has mean

$$\mu_{W^+} = \frac{n(n+1)}{4}$$

and standard deviation

$$\sigma_{W^+} = \sqrt{\frac{n(n+1)(2n+1)}{24}}$$

The **Wilcoxon signed rank test** rejects the hypothesis that there are no systematic differences within pairs when the rank sum  $W^+$  is far from its mean.

### USE YOUR KNOWLEDGE



**15.18 Services provided by top spas.** The readers' poll in *Condé Nast Traveler* magazine that ranked 36 top resort spas and that was described in Exercise 15.1 also reported scores on Diet/Cuisine and on Program/Facilities. Here are the scores for a random sample of 7 spas that ranked in the top 18:

Spa	1	2	3	4	5	6	7
Diet/Cuisine	90.9	92.3	88.6	81.8	85.7	88.9	81.0
Program/Facilities	93.8	92.3	91.4	95.0	89.2	88.2	81.8

Is food, expressed by the Diet/Cuisine score, more important than activities, expressed as the Program/Facilities score, for a top ranking? Formulate this question in terms of null and alternative hypotheses. Then compute the differences and find the value of the Wilcoxon signed rank statistic,  $W^+$ .



**15.19 Scores for lower-ranked spas.** Refer to the previous exercise. Here are the scores for a random sample of 7 spas that ranked between 19 and 36:

Spa	1	2	3	4	5	6	7
Diet/Cuisine	77.3	85.7	84.2	85.3	83.7	84.6	78.5
Program/Facilities	95.7	78.0	87.2	85.3	93.6	76.0	86.3

Answer the questions from the previous exercise for this setting.

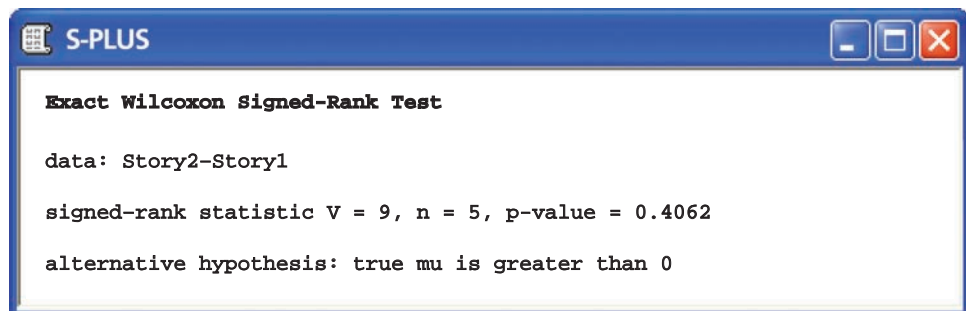
**EXAMPLE**

**15.9 Software output.** In the storytelling study of Example 15.8,  $n = 5$ . If the null hypothesis (no systematic effect of illustrations) is true, the mean of the signed rank statistic is

$$\mu_{W^+} = \frac{n(n+1)}{4} = \frac{(5)(6)}{4} = 7.5$$

Our observed value  $W^+ = 9$  is only slightly larger than this mean. The one-sided  $P$ -value is  $P(W^+ \geq 9)$ .

Figure 15.7 displays the output of two statistical programs. We see from Figure 15.7(a) that the one-sided  $P$ -value for the Wilcoxon signed rank test with  $n = 5$  observations and  $W^+ = 9$  is  $P = 0.4062$ . This result differs from the  $t$  test result  $P = 0.280$ , but both tell us that this very small sample gives no evidence that seeing illustrations improves the storytelling of low-progress readers.



(a)



(b)

**FIGURE 15.7** Output from (a) S-PLUS and (b) SPSS for the storytelling study of Example 15.9. S-PLUS reports the exact  $P$ -value,  $P = 0.4062$ . SPSS uses the Normal approximation without the continuity correction and so gives a less accurate  $P$ -value,  $P = 0.343$  (one-sided).

### The Normal approximation

The distribution of the signed rank statistic when the null hypothesis (no difference) is true becomes approximately Normal as the sample size becomes large. We can then use Normal probability calculations (with the continuity correction) to obtain approximate  $P$ -values for  $W^+$ . Let's see how this works in the storytelling example, even though  $n = 5$  is certainly not a large sample.

**EXAMPLE**

**15.10 The Normal approximation.** For  $n = 5$  observations, we saw in Example 15.9 that  $\mu_{W^+} = 7.5$ . The standard deviation of  $W^+$  under the null hypothesis is

$$\begin{aligned}\sigma_{W^+} &= \sqrt{\frac{n(n+1)(2n+1)}{24}} \\ &= \sqrt{\frac{(5)(6)(11)}{24}} \\ &= \sqrt{13.75} = 3.708\end{aligned}$$

The continuity correction calculates the  $P$ -value  $P(W^+ \geq 9)$  as  $P(W^+ \geq 8.5)$ , treating the value  $W^+ = 9$  as occupying the interval from 8.5 to 9.5. We find the Normal approximation for the  $P$ -value by standardizing and using the standard Normal table:

$$\begin{aligned}P(W^+ \geq 8.5) &= P\left(\frac{W^+ - 7.5}{3.708} \geq \frac{8.5 - 7.5}{3.708}\right) \\ &= P(Z \geq 0.27) \\ &= 0.394\end{aligned}$$

Despite the small sample size, the Normal approximation gives a result quite close to the exact value  $P = 0.4062$ . Figure 15.7(b) shows that the approximation is much less accurate without the continuity correction. *This output reminds us not to trust software unless we know exactly what it does.*

**SPASTOP****SPASNEXT****USE YOUR KNOWLEDGE**

**15.20 Significance test for top-ranked spas.** Refer to Exercise 15.18. Find  $\mu_{W^+}$ ,  $\sigma_{W^+}$ , and the Normal approximation for the  $P$ -value for the Wilcoxon signed rank test.

**15.21 Significance test for lower-ranked spas.** Refer to Exercise 15.19. Find  $\mu_{W^+}$ ,  $\sigma_{W^+}$ , and the Normal approximation for the  $P$ -value for the Wilcoxon signed rank test.

**Ties**

Ties among the absolute differences are handled by assigning average ranks. A tie *within* a pair creates a difference of zero. Because these are neither positive nor negative, the usual procedure simply drops such pairs from the sample. *This amounts to dropping observations that favor the null hypothesis (no difference).* If there are many ties, the test may be biased in favor of the alternative hypothesis. As in the case of the Wilcoxon rank sum, ties complicate finding a  $P$ -value. Most software no longer provides an exact distribution for the signed rank statistic  $W^+$ , and the standard deviation  $\sigma_{W^+}$  must be adjusted for the ties before we can use the Normal approximation. Software will do this. Here is an example.





**EXAMPLE**

**15.11 Golf scores of a women's golf team.** Here are the golf scores of 12 members of a college women's golf team in two rounds of tournament play. (A golf score is the number of strokes required to complete the course, so that low scores are better.)

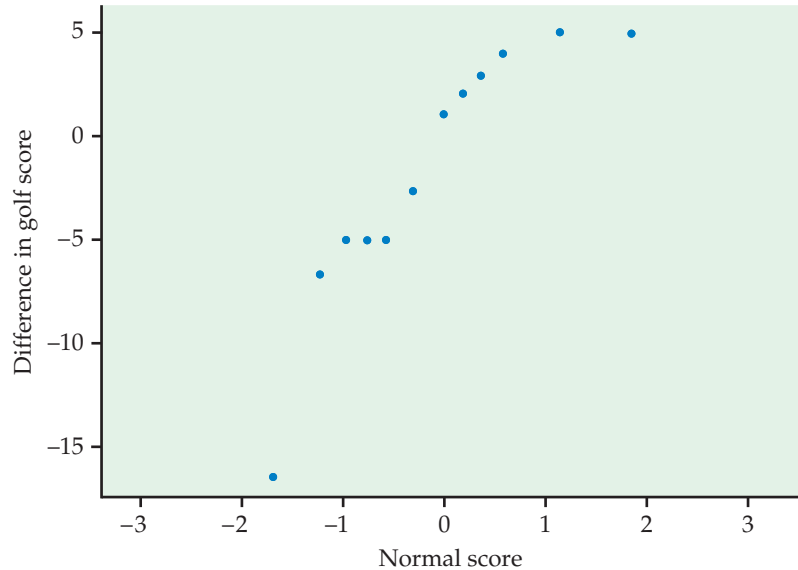
Player	1	2	3	4	5	6	7	8	9	10	11	12
Round 2	94	85	89	89	81	76	107	89	87	91	88	80
Round 1	89	90	87	95	86	81	102	105	83	88	91	79
Difference	5	-5	2	-6	-5	-5	5	-16	4	3	-3	1

Negative differences indicate better (lower) scores on the second round. We see that 6 of the 12 golfers improved their scores. We would like to test the hypotheses that in a large population of collegiate women golfers

$$H_0: \text{Scores have the same distribution in Rounds 1 and 2.}$$

$$H_a: \text{Scores are systematically lower or higher in Round 2.}$$

A Normal quantile plot of the differences (Figure 15.8) shows some irregularity and a low outlier. We will use the Wilcoxon signed rank test.



**FIGURE 15.8** Normal quantile plot of the differences in scores for two rounds of a golf tournament, for Example 15.11.

The absolute values of the differences, with boldface indicating those that were negative, are

5 5 2 **6** 5 5 5 **16** 4 3 3 1

Arrange these in increasing order and assign ranks, keeping track of which values were originally negative. Tied values receive the average of their ranks.

Absolute value	1	2	<b>3</b>	3	4	<b>5</b>	5	<b>5</b>	5	<b>5</b>	<b>6</b>	<b>16</b>
Rank	1	2	<b>3.5</b>	3.5	5	<b>8</b>	8	<b>8</b>	8	<b>8</b>	<b>11</b>	<b>12</b>

The Wilcoxon signed rank statistic is the sum of the ranks of the negative differences. (We could equally well use the sum of the ranks of the positive differences.) Its value is  $W^+ = 50.5$ .

### EXAMPLE

**15.12 Software output.** Here are the two-sided  $P$ -values for the Wilcoxon signed rank test for the golf score data from several statistical programs:

Program	$P$ -value
Minitab	$P = 0.388$
SAS	$P = 0.388$
S-PLUS	$P = 0.384$
SPSS	$P = 0.363$

All lead to the same practical conclusion: these data give no evidence for a systematic change in scores between rounds. However, the  $P$ -values reported differ a bit from program to program. The reason for the variations is that the programs use slightly different versions of the approximate calculations needed when ties are present. The exact result depends on which of these variations the programmer chooses to use.

For these data, the matched pairs  $t$  test gives  $t = 0.9314$  with  $P = 0.3716$ . Once again,  $t$  and  $W^+$  lead to the same conclusion.

## SECTION 15.2 Summary

The **Wilcoxon signed rank test** applies to matched pairs studies. It tests the null hypothesis that there is no systematic difference within pairs against alternatives that assert a systematic difference (either one-sided or two-sided).

The test is based on the **Wilcoxon signed rank statistic**  $W^+$ , which is the sum of the ranks of the positive (or negative) differences when we rank the absolute values of the differences. The **matched pairs  $t$  test** and the **sign test** are alternative tests in this setting.

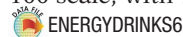
**$P$ -values** for the signed rank test are based on the sampling distribution of  $W^+$  when the null hypothesis is true. You can find  $P$ -values from special tables, software, or a Normal approximation (with continuity correction).

## SECTION 15.2 Exercises

For Exercises 15.18 and 15.19, see page 15-18; and for Exercises 15.20 and 15.21, see page 15-20.


Statistical software is very helpful in doing these exercises. If you do not have access to software, base your work on the Normal approximation with continuity correction (page 15-20).

**15.22 Comparison of two energy drinks.** Consider the following study to compare two popular energy drinks. For each subject, a coin was flipped to determine which drink to rate first. Each drink was rated on a 0 to 100 scale, with 100 being the highest rating.




Drink	Subject					
	1	2	3	4	5	6
A	43	83	66	87	78	67
B	45	78	64	79	71	62

- (a) Inspect the data. Is there a tendency for these subjects to prefer one of the two energy drinks?
- (b) Use the matched pairs  $t$  test of Chapter 7 (page 424) to compare the two drinks.
- (c) Use the Wilcoxon signed rank test to compare the two drinks.
- (d) Write a summary of your results and explain why the two tests give different conclusions.

**15.23 Comparison of two energy drinks with an additional subject.** Refer to the previous exercise. Let's suppose that there is an additional subject who expresses a strong preference for energy drink "A." Here is the new data set:  ENERGYDRINKS7


Drink	Subject						
	1	2	3	4	5	6	7
A	43	83	66	87	78	67	90
B	45	78	64	79	71	62	60

Answer the questions given in the previous exercise. Write a summary comparing this exercise with the previous one. Include a discussion of what you have learned regarding the choice of the  $t$  test versus the Wilcoxon signed rank test for different sets of data.

**15.24 Carbon dioxide and plant growth.** The concentration of carbon dioxide ( $\text{CO}_2$ ) in the atmosphere is increasing rapidly due to our use of fossil fuels. Because plants use  $\text{CO}_2$  to fuel photosynthesis, more  $\text{CO}_2$  may cause trees and other plants to grow faster. An elaborate apparatus allows researchers to pipe extra  $\text{CO}_2$  to a 30-meter circle of forest. They set up three pairs of circles in different parts of a forest in North Carolina. One of each pair received extra  $\text{CO}_2$  for an entire growing season, and the other received ambient air. The response variable is the average growth in base area for trees in a circle, as a fraction of the starting area. Here are the data for one growing season:<sup>12</sup>  CO2PLANTS


Pair	Control	Treatment
1	0.06528	0.08150
2	0.05232	0.06334
3	0.04329	0.05936

- (a) Summarize the data. Does it appear that growth was faster in the treated plots?
- (b) The researchers used a matched pairs  $t$  test to see if the data give good evidence of faster growth in the treated plots. State hypotheses, carry out the test, and state your conclusion.
- (c) The sample is so small that we cannot assess Normality. To be safe, we might use the Wilcoxon signed rank test. Carry out this test and report your result.
- (d) The tests lead to very different conclusions. The primary reason is the lack of power of rank tests for very small samples. Explain to someone who knows no statistics what this means.

**15.25 Heart rate and exercise.** A student project asked subjects to step up and down for three minutes and measured their heart rates before and after the exercise. Here are data for five subjects and two treatments: stepping at a low rate (14 steps per minute) and at a medium rate (21 steps per minute). For each subject, we give the resting heart rate (beats per minute) and the heart rate at the end of the exercise.<sup>13</sup>  HEARTEXERCISE

Subject	Low Rate		Medium Rate	
	Resting	Final	Resting	Final
1	60	75	63	84
2	90	99	69	93
3	87	93	81	96
4	78	87	75	90
5	84	84	90	108

Does exercise at the low rate raise heart rate significantly? State hypotheses in terms of the median increase in heart rate and apply the Wilcoxon signed rank test. What do you conclude?

**15.26 Compare exercise at a medium rate with exercise at a low rate.** Do the data from the previous exercise give good reason to think that stepping at the medium rate increases heart rates more than stepping at the low rate?  HEARTEXERCISE


- (a) State hypotheses in terms of comparing the median increases for the two treatments. What is the proper rank test for these hypotheses?
- (b) Carry out your test and state a conclusion.

**15.27 The full moon and behavior.** Can the full moon influence behavior? A study observed 15 nursing-home patients with dementia. The number of incidents of aggressive behavior was recorded each day for 12 weeks. Call a day a "moon day" if it is the day of a full moon or the day before or after a full moon. Here are the average

numbers of aggressive incidents for moon days and other days for each subject:<sup>14</sup>  MOONBEHAVIOR

Patient	Moon days	Other days
1	3.33	0.27
2	3.67	0.59
3	2.67	0.32
4	3.33	0.19
5	3.33	1.26
6	3.67	0.11
7	4.67	0.30
8	2.67	0.40
9	6.00	1.59
10	4.33	0.60
11	3.33	0.65
12	0.67	0.69
13	1.33	1.26
14	0.33	0.23
15	2.00	0.38


The matched pairs *t* test (Example 7.7, page 414) gives  $P < 0.000015$  and a permutation test (Example 16.14, page 16-48) gives  $P = 0.0001$ . Does the Wilcoxon signed rank test, based on ranks rather than means, agree that there is strong evidence that there are more aggressive incidents on moon days?

**15.28 A summer language institute for teachers.** A matched pairs study of the effect of a summer language institute on the ability of teachers to comprehend spoken French had these improvements in scores between the pretest and the posttest for 20 teachers:  SUMMERLANGUAGE

2	0	6	6	3	3	2	3	-6	6
6	6	3	0	1	1	0	2	3	3

(Exercise 7.41, page 431, applies the *t* test to these data; Exercise 16.59, page 16-52, applies a permutation test based on the means.) Show the assignment of ranks and the calculation of the signed rank statistic  $W^+$  for these data. Remember that zeros are dropped from the data before ranking, so that  $n$  is the number of nonzero differences within pairs.

**15.29 Radon detectors.** How accurate are radon detectors of a type sold to home owners? To answer this question, university researchers placed 12 detectors in a chamber that exposed them to 105 picocuries per liter (pCi/l) of radon.<sup>15</sup> The detector readings are as follows:

 RADONDETECTORS

91.9	97.8	111.4	122.3	105.4	95.0
103.8	99.6	96.6	119.3	104.8	101.7

We wonder if the median reading differs significantly from the true value 105.


(a) Graph the data, and comment on skewness and outliers. A rank test is appropriate.

(b) We would like to test hypotheses about the median reading from home radon detectors:

$$H_0 : \text{median} = 105$$


$$H_a : \text{median} \neq 105$$

To do this, apply the Wilcoxon signed rank statistic to the differences between the observations and 105. (This is the one-sample version of the test.) What do you conclude?

**15.30 Vitamin C in wheat-soy blend.** The U.S. Agency for International Development provides large quantities of wheat-soy blend (WSB) for development programs and emergency relief in countries throughout the world. One study collected data on the vitamin C content of 27 bags of WSB at the factory and five months later in Haiti.<sup>16</sup> Here are the data:  WSBVITC

Sample	1	2	3	4	5
Before	73	79	86	88	78
After	20	27	29	36	17

We want to know if vitamin C has been lost during transportation and storage. Describe what the data show about this question. Then use a rank test to see whether there has been a significant loss.

**15.31 Weight gains with an extra 1000 calories per day.** Exercise 7.32 (page 428) presents these data on the weight gains (in kilograms) of adults who were fed an extra 1000 calories per day for 8 weeks:<sup>17</sup>  WEIGHT1000

Subject	Weight	
	before	after
1	55.7	61.7
2	54.9	58.8
3	59.6	66.0
4	62.3	66.2
5	74.2	79.0
6	75.6	82.3
7	70.7	74.3
8	53.3	59.3
9	73.3	79.1
10	63.4	66.0
11	68.1	73.4
12	73.7	76.9
13	91.7	93.1
14	55.9	63.0
15	61.7	68.2
16	57.8	60.3

(a) Use a rank test to test the null hypothesis that the median weight gain is 16 pounds, as theory suggests. What do you conclude?

(b) If your software allows, give a 95% confidence interval for the median weight gain in the population.

## 15.3 The Kruskal-Wallis Test\*

We have now considered alternatives to the two-sample  $t$  and matched pairs tests for comparing the magnitude of responses to two treatments. To compare more than two treatments, we use one-way analysis of variance (ANOVA) if the distributions of the responses to each treatment are at least roughly Normal and have similar spreads. What can we do when these distribution requirements are violated?

### EXAMPLE

**15.13 Weeds and corn yield.** Lamb's-quarter is a common weed that interferes with the growth of corn. A researcher planted corn at the same rate in 16 small plots of ground and then randomly assigned the plots to four groups. He weeded the plots by hand to allow a fixed number of lamb's-quarter plants to grow in each meter of corn row. These numbers were 0, 1, 3, and 9 in the four groups of plots. No other weeds were allowed to grow, and all plots received identical treatment except for the weeds. Here are the yields of corn (bushels per acre) in each of the plots:<sup>18</sup>



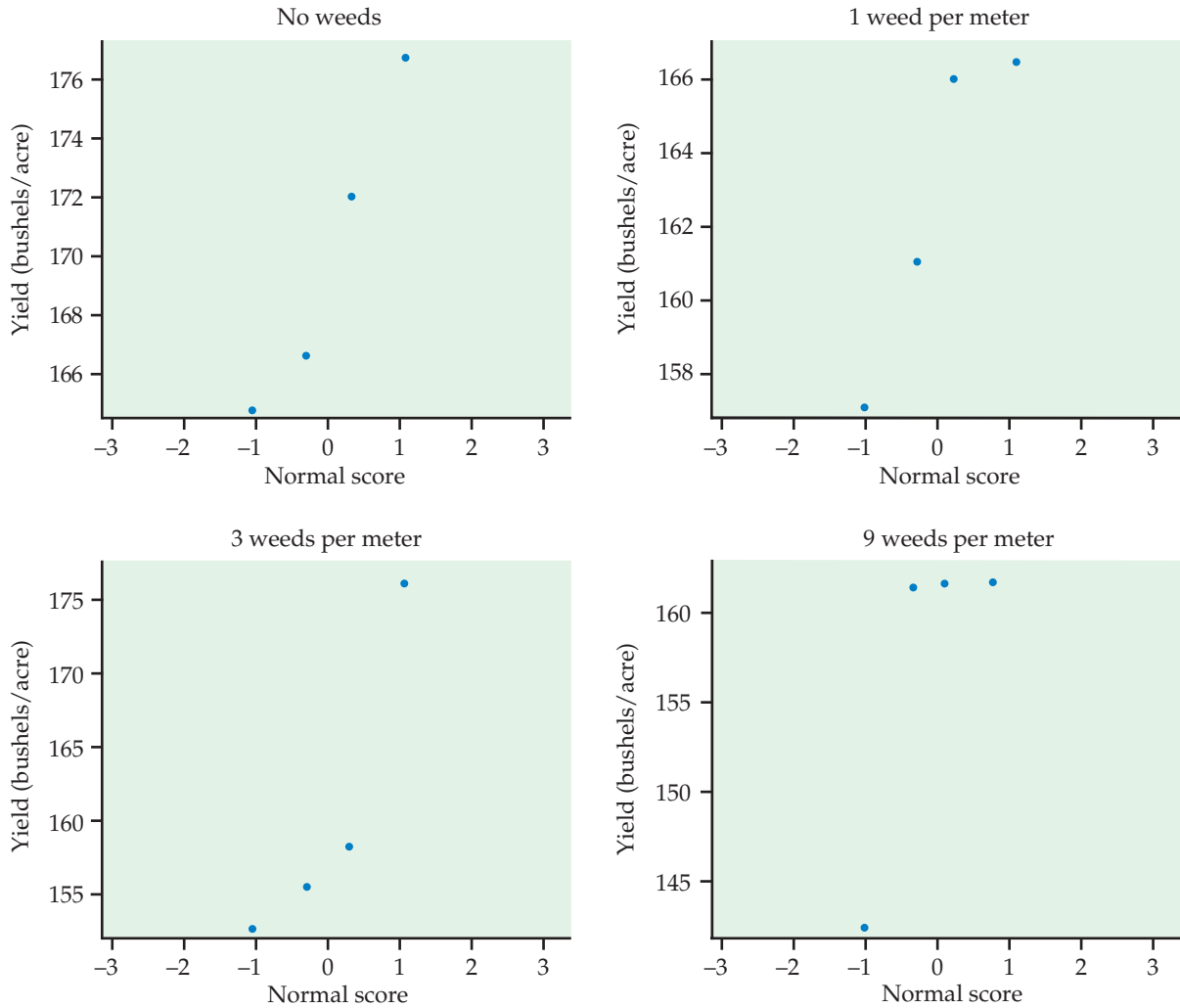
Weeds per meter	Corn yield	Weeds per meter	Corn yield	Weeds per meter	Corn yield	Weeds per meter	Corn yield
0	166.7	1	166.2	3	158.6	9	162.8
0	172.2	1	157.3	3	176.4	9	142.4
0	165.0	1	166.7	3	153.1	9	162.7
0	176.9	1	161.1	3	156.0	9	162.4

The summary statistics are

Weeds	$n$	Mean	Std. dev.
0	4	170.200	5.422
1	4	162.825	4.469
3	4	161.025	10.493
9	4	157.575	10.118

The sample standard deviations do not satisfy our rule of thumb that for safe use of ANOVA the largest should not exceed twice the smallest. Normal quantile plots (Figure 15.9) show that outliers are present in the yields for three and nine weeds per meter. These are the correct yields for their plots, so we have no justification for removing them. We may want to use a rank test.

\*Because this test is an alternative to the one-way analysis of variance  $F$  test, you should first read Chapter 12.



**FIGURE 15.9** Normal quantile plots for the corn yields in the four treatment groups in Example 15.13.

### Hypotheses and assumptions

The ANOVA  $F$  test concerns the means of the several populations represented by our samples. In Example 15.13, the ANOVA hypotheses are

$$H_0: \mu_0 = \mu_1 = \mu_3 = \mu_9$$

$$H_a: \text{not all four means are equal}$$

Here,  $\mu_0$  is the mean yield in the population of all corn planted under the conditions of the experiment with no weeds present. The data should consist of four independent random samples from the four populations, all Normally distributed with the same standard deviation.

The *Kruskal-Wallis test* is a rank test that can replace the ANOVA  $F$  test. The assumption about data production (independent random samples from each population) remains important, but we can relax the Normality assumption.

We assume only that the response has a continuous distribution in each population. The hypotheses tested in our example are

$H_0$ : Yields have the same distribution in all groups.

$H_a$ : Yields are systematically higher in some groups than in others.

If all of the population distributions have the same shape (Normal or not), these hypotheses take a simpler form. The null hypothesis is that all four populations have the same *median* yield. The alternative hypothesis is that not all four median yields are equal.

### The Kruskal-Wallis test

Recall the analysis of variance idea: we write the total observed variation in the responses as the sum of two parts, one measuring variation among the groups (sum of squares for groups, SSG) and one measuring variation among individual observations within the same group (sum of squares for error, SSE). The ANOVA  $F$  test rejects the null hypothesis that the mean responses are equal in all groups if SSG is large relative to SSE.

The idea of the Kruskal-Wallis rank test is to rank all the responses from all groups together and then apply one-way ANOVA to the ranks rather than to the original observations. If there are  $N$  observations in all, the ranks are always the whole numbers from 1 to  $N$ . The total sum of squares for the ranks is therefore a fixed number no matter what the data are. So we do not need to look at both SSG and SSE. Although it isn't obvious without some unpleasant algebra, the Kruskal-Wallis test statistic is essentially just SSG for the ranks. We give the formula, but you should rely on software to do the arithmetic. When SSG is large, that is evidence that the groups differ.

#### THE KRUSKAL-WALLIS TEST

Draw independent SRSs of sizes  $n_1, n_2, \dots, n_I$  from  $I$  populations. There are  $N$  observations in all. Rank all  $N$  observations and let  $R_i$  be the sum of the ranks for the  $i$ th sample. The **Kruskal-Wallis statistic** is

$$H = \frac{12}{N(N+1)} \sum \frac{R_i^2}{n_i} - 3(N+1)$$

When the sample sizes  $n_i$  are large and all  $I$  populations have the same continuous distribution,  $H$  has approximately the chi-square distribution with  $I - 1$  degrees of freedom.

The **Kruskal-Wallis test** rejects the null hypothesis that all populations have the same distribution when  $H$  is large.

We now see that, like the Wilcoxon rank sum statistic, the Kruskal-Wallis statistic is based on the sums of the ranks for the groups we are comparing. The more different these sums are, the stronger is the evidence that responses are systematically larger in some groups than in others.

The exact distribution of the Kruskal-Wallis statistic  $H$  under the null hypothesis depends on all the sample sizes  $n_1$  to  $n_I$ , so tables are awkward.

The calculation of the exact distribution is so time-consuming for all but the smallest problems that even most statistical software uses the chi-square approximation to obtain  $P$ -values. As usual, there is no usable exact distribution when there are ties among the responses. We again assign average ranks to tied observations.



**EXAMPLE**

**15.14 Perform the significance test.** In Example 15.13, there are  $I = 4$  populations and  $N = 16$  observations. The sample sizes are equal,  $n_i = 4$ . The 16 observations arranged in increasing order, with their ranks, are

Yield	142.4	153.1	156.0	157.3	158.6	161.1	162.4	162.7
Rank	1	2	3	4	5	6	7	8
Yield	162.8	165.0	166.2	166.7	166.7	172.2	176.4	176.9
Rank	9	10	11	12.5	12.5	14	15	16

There is one pair of tied observations. The ranks for each of the four treatments are

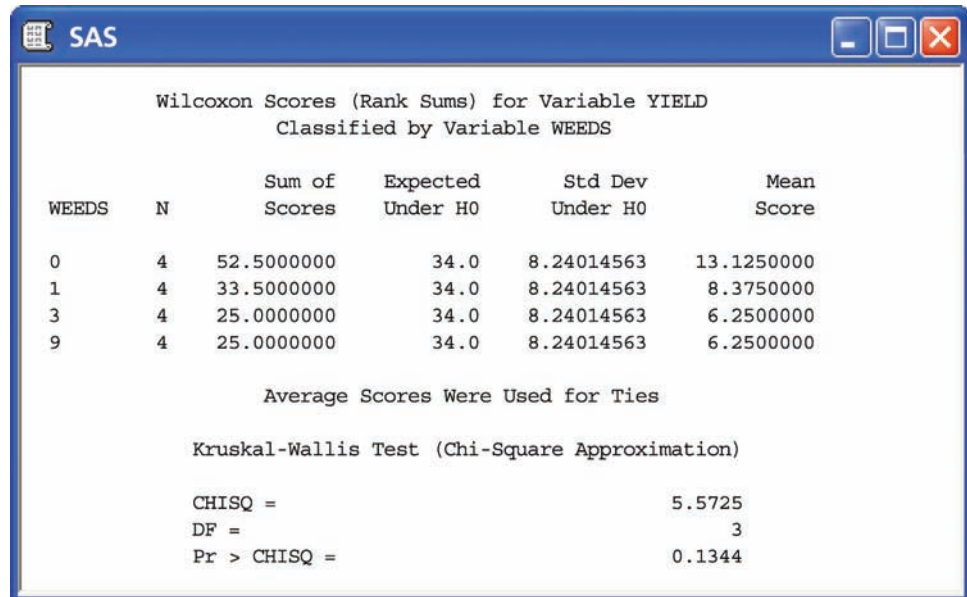
Weeds	Ranks				Rank sums
0	10	12.5	14	16	52.5
1	4	6	11	12.5	33.5
3	2	3	5	15	25.0
9	1	7	8	9	25.0

The Kruskal-Wallis statistic is therefore

$$\begin{aligned}
 H &= \frac{12}{N(N+1)} \sum \frac{R_i^2}{n_i} - 3(N+1) \\
 &= \frac{12}{(16)(17)} \left( \frac{52.5^2}{4} + \frac{33.5^2}{4} + \frac{25^2}{4} + \frac{25^2}{4} \right) - (3)(17) \\
 &= \frac{12}{272} (1282.125) - 51 \\
 &= 5.56
 \end{aligned}$$

Referring to the table of chi-square critical points (Table F) with  $df = 3$ , we find that the  $P$ -value lies in the interval  $0.10 < P < 0.15$ . This small experiment suggests that more weeds decrease yield but does not provide convincing evidence that weeds have an effect.

Figure 15.10 displays the output from the SAS statistical software, which gives the results  $H = 5.5725$  and  $P = 0.1344$ . The software makes a small



**FIGURE 15.10** Output from SAS for the Kruskal-Wallis test applied to the data in Example 15.14. SAS uses the chi-square approximation to obtain a  $P$ -value.

adjustment for the presence of ties that accounts for the slightly larger value of  $H$ . The adjustment makes the chi-square approximation more accurate. It would be important if there were many ties.

As an option, SAS will calculate the exact  $P$ -value for the Kruskal-Wallis test. The result for Example 15.14 is  $P = 0.1299$ . This result required more than an hour of computing time. Fortunately, the chi-square approximation is quite accurate. The ordinary ANOVA  $F$  test gives  $F = 1.73$  with  $P = 0.2130$ . Although the practical conclusion is the same, ANOVA and Kruskal-Wallis do not agree closely in this example. The rank test is more reliable for these small samples with outliers.

## SECTION 15.3 Summary

The **Kruskal-Wallis test** compares several populations on the basis of independent random samples from each population. This is the **one-way analysis of variance** setting.


The null hypothesis for the Kruskal-Wallis test is that the distribution of the response variable is the same in all the populations. The alternative hypothesis is that responses are systematically larger in some populations than in others.

The **Kruskal-Wallis statistic  $H$**  can be viewed in two ways. It is essentially the result of applying one-way ANOVA to the ranks of the observations. It is also a comparison of the sums of the ranks for the several samples.


When the sample sizes are not too small and the null hypothesis is true,  $H$  for comparing  $I$  populations has approximately the chi-square distribution with  $I - 1$  degrees of freedom. We use this approximate distribution to obtain  $P$ -values.

**SECTION 15.3 Exercises**

Statistical software is needed to do these exercises without unpleasant hand calculations. If you do not have access to software, find the Kruskal-Wallis statistic  $H$  by hand and use the chi-square table to get approximate  $P$ -values.


**15.32 Do poets die young?** In Exercise 12.38 you analyzed the age at death for female writers. They were classified as novelists, poets, and nonfiction writers. The data are given in Table 12.1 (page 662).  POETS

- (a) Use the Kruskal-Wallis test to compare the three groups of female writers.
- (b) Compare these results with what you find using the ANOVA  $F$  statistic.

**15.33 Do isoflavones increase bone mineral density?** In Exercise 12.39 (page 662) you investigated the effects of isoflavones from Kudzu on bone mineral density (BMD). The experiment randomized rats to three diets: control, low isoflavones, and high isoflavones. Here are the data:  BMD

Treatment	BMD (g/cm <sup>2</sup> )														
Control	0.228	0.207	0.234	0.220	0.217	0.228	0.209	0.221	0.204	0.220	0.203	0.219	0.218	0.245	0.210
Low dose	0.211	0.220	0.211	0.233	0.219	0.233	0.226	0.228	0.216	0.225	0.200	0.208	0.198	0.208	0.203
High dose	0.250	0.237	0.217	0.206	0.247	0.228	0.245	0.232	0.267	0.261	0.221	0.219	0.232	0.209	0.255


- (a) Use the Kruskal-Wallis test to compare the three diets.
- (b) How do these results compare with what you find using the ANOVA  $F$  statistic?

**15.34 Vitamins in bread.** Does bread lose its vitamins when stored? Here are data on the vitamin C content (milligrams per 100 grams of flour) in bread baked from the same recipe and stored for 1, 3, 5, or 7 days.<sup>19</sup> The 10 observations are from 10 different loaves of bread.  BREAD

Condition	Vitamin C (mg/100 g)	
Immediately after baking	47.62	49.79
One day after baking	40.45	43.46
Three days after baking	21.25	22.34
Five days after baking	13.18	11.65
Seven days after baking	8.51	8.13


The loss of vitamin C over time is clear, but with only 2 loaves of bread for each storage time we wonder if the differences among the groups are significant.

- (a) Use the Kruskal-Wallis test to assess significance and then write a brief summary of what the data show.
- (b) Because there are only 2 observations per group, we suspect that the common chi-square approximation to the distribution of the Kruskal-Wallis statistic may not be accurate. The exact  $P$ -value (from the SAS software) is  $P = 0.0011$ . Compare this with your  $P$ -value from part (a). Is the difference large enough to affect your conclusion?

**15.35 Jumping and strong bones.** Many studies suggest that exercise causes bones to get stronger. One study examined the effect of jumping on the bone density of growing rats. Ten rats were assigned to each of three treatments: a 60-centimeter “high jump,” a 30-centimeter “low jump,” and a control group with no jumping. Here are the bone densities (in milligrams per cubic centimeter) after eight weeks of 10 jumps per day:<sup>20</sup>  JUMPING

Group	Bone density (mg/cm <sup>3</sup> )									
Control	611	621	614	593	593	653	600	554	603	569
Low jump	635	605	638	594	599	632	631	588	607	596
High jump	650	622	626	626	631	622	643	674	643	650


- (a) The study was a randomized comparative experiment. Outline the design of this experiment.
- (b) Make side-by-side stemplots for the three groups, with the stems lined up for easy comparison. The distributions are a bit irregular but not strongly non-Normal. We would usually use analysis of variance to assess the significance of the difference in group means.
- (c) Do the Kruskal-Wallis test. Explain the distinction between the hypotheses tested by Kruskal-Wallis and ANOVA.
- (d) Write a brief statement of your findings. Include a numerical comparison of the groups as well as your test result.

**15.36 Detecting insects in farm fields.** To detect the presence of harmful insects in farm fields, we can put up boards covered with a sticky material and examine the insects trapped on the boards. Which colors attract insects best? Experimenters placed six boards of each of four colors at random locations in a field of oats and measured the number of cereal leaf beetles trapped. Here are the data:<sup>21</sup>  INSECTS


Color	Insects trapped					
Lemon yellow	45	59	48	46	38	47
White	21	12	14	17	13	17
Green	37	32	15	25	39	41
Blue	16	11	20	21	14	7

Because the samples are small, we will apply a nonparametric test.

- (a) What hypotheses does ANOVA test? What hypotheses does Kruskal-Wallis test?
- (b) Find the median number of beetles trapped by boards of each color. Which colors appear more effective? Use the Kruskal-Wallis test to see if there are significant differences among the colors. What do you conclude?


**15.37 Do the calculations by hand.** Exercise 15.36 gives data on the counts of insects attracted by boards of four different colors. Carry out the Kruskal-Wallis test by hand, following these steps.  INSECTS

- (a) What are  $I$ , the  $n_i$ , and  $N$ ?
- (b) Arrange the counts in order and assign ranks. Be careful about ties. Find the sum of the ranks  $R_i$  for each color.
- (c) Calculate the Kruskal-Wallis statistic  $H$ . How many degrees of freedom should you use for the chi-square approximation to its null distribution? Use the chi-square table to give an approximate  $P$ -value.

**15.38 Logging in Borneo.** In Exercise 15.15 you compared the number of tree species in plots of land in a tropical rainforest that had never been logged with similar plots nearby that had been logged eight years earlier. The researchers also counted species in plots that had been logged just one year earlier. Here are the counts of species:<sup>22</sup>  BORNEO3

Plot type	Species count					
Unlogged	22	18	22	20	15	21
	13	13	19	13	19	15
Logged 1 year ago	11	11	14	7	18	15
	15	12	13	2	15	8
Logged 8 years ago	17	4	18	14	18	15
	15	10	12			

- (a) Use side-by-side stemplots to compare the distributions of number of species per plot for the three groups of plots. Are there features that might prevent use of ANOVA? Also give the median number of species per plot in the three groups.
- (b) Use the Kruskal-Wallis test to compare the distributions of species counts. State hypotheses, the test statistic and its  $P$ -value, and your conclusions.


**15.39 Heart disease and smoking.** In a study of heart disease in male federal employees, researchers classified 356 volunteer subjects according to their socioeconomic status (SES) and their smoking habits. There were three categories of SES: high, middle, and low. Individuals were asked whether they were current smokers, former smokers, or had never smoked. Here are the data, as a two-way table of counts:<sup>23</sup>  SMOKINGSES

SES	Never (1)	Former (2)	Current (3)
High	68	92	51
Middle	9	21	22
Low	22	28	43

Smoking behavior is stored numerically as 1, 2, or 3 using the codes given in the column headings above.

- (a) Higher-SES people in the United States smoke less as a group than lower-SES people. Do these data show a relationship of this kind? Give percents that back your statements.
- (b) Apply the chi-square test to see if there is a significant relationship between SES and smoking behavior.
- (c) The chi-square test ignores the ordering of the responses. Use the Kruskal-Wallis test (with many ties) to test the hypothesis that some SES classes smoke systematically more than others.

## CHAPTER 15 Exercises

**15.40 Time spent studying.** In Exercise 1.41 (page 26) you compared the time spent studying by men and women. The students in a large first-year college class were asked how many minutes they studied on a typical weeknight. Here are the responses of random samples of 30 women and 30 men from the class:  STUDYTIME

Women					Men				
180	120	180	360	240	90	120	30	90	200
120	180	120	240	170	90	45	30	120	75
150	120	180	180	150	150	120	60	240	300
200	150	180	150	180	240	60	120	60	30
120	60	120	180	180	30	230	120	95	150
90	240	180	115	120	0	200	120	120	180

- (a) Summarize the data numerically and graphically.
- (b) Use the Wilcoxon rank sum test to compare the men and women. Write a short summary of your results.
- (c) Use a two-sample  $t$  test to compare the men and women. Write a short summary of your results.
- (d) Which procedure is more appropriate for these data? Give reasons for your answer.

**15.41 Response times for telephone repair calls.** A study examined on the time required for the telephone company Verizon to respond to repair calls from its own

customers and from customers of a CLEC, another phone company that pays Verizon to use its local lines. Here are the data, which are rounded to the nearest hour:



Verizon											
1	1	1	1	2	2	1	1	1	1	2	2
1	1	1	1	2	2	1	1	1	1	2	3
1	1	1	1	2	3	1	1	1	1	2	3
1	1	1	1	2	3	1	1	1	1	2	3
1	1	1	1	2	3	1	1	1	1	2	4
1	1	1	1	2	5	1	1	1	1	2	5
1	1	1	1	2	6	1	1	1	1	2	8
1	1	1	1	2	15	1	1	1	2	2	

CLEC					
1	1	5	5	5	5

- (a) Does Verizon appear to give CLEC customers the same level of service as its own customers? Compare the data using graphs and descriptive measures and express your opinion.
- (b) We would like to see if times are significantly longer for CLEC customers than for Verizon customers. Why would you hesitate to use a  $t$  test for this purpose? Carry out a rank test. What can you conclude?

**15.42 Selling prices of three- and four-bedroom homes.** Exercise 7.141 (page 471) reports data on the selling prices of 14 four-bedroom houses and 23 three-bedroom houses in West Lafayette, Indiana. We wonder if there is a difference between the average prices of three- and four-bedroom houses in this community.



- (a) Make a Normal quantile plot of the prices of three-bedroom houses. What kind of deviation from Normality do you see?
- (b) The  $t$  tests are quite robust. State the hypotheses for the proper  $t$  test, carry out the test, and present your results, including appropriate data summaries.
- (c) Carry out a nonparametric test. Once more state the hypotheses tested and present your results for both the test and the data summaries that should go with them.

**15.43 Plants and hummingbirds.** Different varieties of the tropical flower *Heliconia* are fertilized by different species of hummingbirds. Over time, the lengths of the flowers and the form of the hummingbirds' beaks have evolved to match each other. Here are data on the lengths in millimeters of three varieties of these flowers on the island of Dominica:<sup>24</sup>



<i>H. bihai</i>					
47.12	46.75	46.81	47.12	46.67	47.43
46.44	46.64	48.07	48.34	48.15	50.26
50.12	46.34	46.94	48.36		

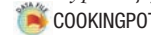
<i>H. caribaea red</i>					
41.90	42.01	41.93	43.09	41.47	41.69
39.78	40.57	39.63	42.18	40.66	37.87
39.16	37.40	38.20	38.07	38.10	37.97
38.79	38.23	38.87	37.78	38.01	

<i>H. caribaea yellow</i>					
36.78	37.02	36.52	36.11	36.03	35.45
38.13	37.10	35.17	36.82	36.66	35.68
36.03	34.57	34.63			

Do a complete analysis that includes description of the data and a rank test for the significance of the differences in lengths among the three species.

*Iron-deficiency anemia is the most common form of malnutrition in developing countries. Does the type of cooking pot affect the iron content of food? We have data from a study in Ethiopia that measured the iron content (milligrams per 100 grams of food) for three types of food cooked in each of three types of pots:<sup>25</sup>*



Type of pot	Iron Content			
	Meat			
Aluminum	1.77	2.36	1.96	2.14
Clay	2.27	1.28	2.48	2.68
Iron	5.27	5.17	4.06	4.22
	Legumes			
Aluminum	2.40	2.17	2.41	2.34
Clay	2.41	2.43	2.57	2.48
Iron	3.69	3.43	3.84	3.72
	Vegetables			
Aluminum	1.03	1.53	1.07	1.30
Clay	1.55	0.79	1.68	1.82
Iron	2.45	2.99	2.80	2.92

Exercises 15.44 to 15.46 use these data.


**15.44 Cooking vegetables in different pots.** Does the vegetable dish vary in iron content when cooked in aluminum, clay, and iron pots?



- (a) What do the data appear to show? Check the conditions for one-way ANOVA. Which requirements are a bit dubious in this setting?
- (b) Instead of ANOVA, do a rank test. Summarize your conclusions about the effect of pot material on the iron content of the vegetable dish.

**15.45 Cooking meat and legumes in aluminum and clay pots.** There appears to be little difference between the iron content of food cooked in aluminum pots and food cooked in clay pots. Is there a significant difference between the iron content of meat cooked in aluminum and clay? Is the difference between aluminum and clay significant for legumes? Use rank tests. 🍳 COOKINGPOT

**15.46 Iron in food cooked in iron pots.** The data show that food cooked in iron pots has the highest iron content. They also suggest that the three types of food differ in iron content. Is there significant evidence that the three types of food differ in iron content when all are cooked in iron pots? 🍳 COOKINGPOT


**15.47**  **Multiple comparisons for plants and hummingbirds.** As in ANOVA, we often want to carry out a **multiple-comparisons** procedure following a Kruskal-Wallis test to tell us *which* groups differ significantly.<sup>26</sup> Here is a simple method: If we carry out  $k$  tests at fixed significance level  $0.05/k$ , the probability of *any* false rejection among the  $k$  tests is always no greater than 0.05. That is, to get overall significance level 0.05 for

all of  $k$  comparisons, do each individual comparison at the  $0.05/k$  level. In Exercise 15.43 you found a significant difference among the lengths of three varieties of the flower *Heliconia*. Now we will explore multiple comparisons. 🍷 HUMMINGBIRDS

(a) Write down all the pairwise comparisons we can make, for example, *bihai* versus *caribaea* red. There are three possible pairwise comparisons.

(b) Carry out three Wilcoxon rank sum tests, one for each of the three pairs of flower varieties. What are the three two-sided  $P$ -values?

(c) For purposes of multiple comparisons, any of these three tests is significant if its  $P$ -value is no greater than  $0.05/3 = 0.0167$ . Which pairs differ significantly at the overall 0.05 level?

**15.48**  **Multiple comparisons for cooking pots.** The previous exercise outlines how to use the Wilcoxon rank sum test several times for multiple comparisons with overall significance level 0.05 for all comparisons together. Apply this procedure to the data used in each of Exercises 15.44 to 15.46. 🍳 COOKINGPOT

## Chapter 15 Notes

1. Data provided by Sam Phillips, Purdue University.
2. From the April 2007 issue of *Condé Nast Traveler* magazine.
3. For purists, here is the precise definition:  $X_1$  is *stochastically larger* than  $X_2$  if

$$P(X_1 > a) \geq P(X_2 > a)$$

for all  $a$ , with strict inequality for at least one  $a$ . The Wilcoxon rank sum test is effective against this alternative in the sense that the power of the test approaches 1 (that is, the test becomes more certain to reject the null hypothesis) as the number of observations increases.

4. Based on a Pew Research Center Report entitled “Take this job and love it,” by Rich Morin, September 17, 2009. See [pewsocialtrends.org/pubs](http://pewsocialtrends.org/pubs)
5. From Matthias R. Mehl et al., “Are women really more talkative than men?” *Science*, 317 (5834), (2007), p. 82. The raw data were provided by Matthias Mehl.
6. Data provided by Susan Stadler, Purdue University.
7. Data provided by Warren Page, New York City Technical College, from a study done by John Hudesman.
8. Data provided by Charles Cannon, Duke University. The study report is C. H. Cannon, D. R. Peart, and M. Leighton, “Tree species diversity in commercially logged Bornean rainforest,” *Science*, 281 (1998), pp. 1366–1367.
9. This example is adapted from Maribeth C. Schmitt, “The effects of an elaborated directed reading activity on the metacomprehension skills of third graders,” PhD dissertation, Purdue University, 1987.
10. William D. Darley, “Store-choice behavior for pre-owned merchandise,” *Journal of Business Research*, 27 (1993), pp. 17–31.
11. See Note 5.
12. Data for 1998 provided by Jason Hamilton, University of Illinois. The study report is Evan H. DeLucia et al., “Net primary production of a forest ecosystem with experimental CO<sub>2</sub> enhancement,” *Science*, 284 (1999), pp. 1177–1179.

13. Simplified from the EESEE story “Stepping Up Your Heart Rate,” on the course Web site.
14. These data were collected as part of a larger study of dementia patients conducted by Nancy Edwards, School of Nursing, and Alan Beck, School of Veterinary Medicine, Purdue University.
15. Data provided by Diana Schellenberg, Purdue University School of Health Sciences.
16. These data are from “Results report on the vitamin C pilot program,” prepared by SUSTAIN (Sharing United States Technology to Aid in the Improvement of Nutrition) for the U.S. Agency for International Development. The report was used by the Committee on International Nutrition of the National Academy of Sciences/Institute of Medicine to make recommendations on whether or not the vitamin C content of food commodities used in U.S. food aid programs should be increased. The program was directed by Peter Ranum and Françoise Chomé.
17. James A. Levine et al., “Role of nonexercise activity thermogenesis in resistance to fat gain in humans,” *Science*, 283 (1999), pp. 212–214. Data for this study are available from the *Science* Web site, [sciencemag.org](http://sciencemag.org)
18. See Note 1.
19. Data provided by Helen Park. See H. Park et al., “Fortifying bread with each of three antioxidants,” *Cereal Chemistry*, 74 (1997), pp. 202–206.
20. Data provided by Jo Welch, Purdue University Department of Foods and Nutrition.
21. Modified from M. C. Wilson and R. E. Shade, “Relative attractiveness of various luminescent colors to the cereal leaf beetle and the meadow spittlebug,” *Journal of Economic Entomology*, 60 (1967), pp. 578–580.
22. See Note 8.
23. Ray H. Rosenman et al., “A 4-year prospective study of the relationship of different habitual vocational physical activity to risk and incidence of ischemic heart disease in volunteer male federal employees,” in P. Milvey (ed.), *The Marathon: Physiological, Medical, Epidemiological and Psychological Studies*, New York Academy of Sciences, 301 (1977), pp. 627–641.
24. We thank Ethan J. Temeles of Amherst College for providing the data. His work is described in Ethan J. Temeles and W. John Kress, “Adaptation in a plant-hummingbird association,” *Science*, 300 (2003), pp. 630–633.
25. Based on A. A. Adish et al., “Effect of consumption of food cooked in iron pots on iron status and growth of young children: A randomised trial,” *The Lancet*, 353 (1999), pp. 712–716.
26. For more details on multiple comparisons, see M. Hollander and D. A. Wolfe, *Nonparametric Statistical Methods*, 2nd ed., Wiley, 1999. This book is a useful reference on applied aspects of nonparametric inference in general.