

CHAPTER 14

Nonparametric Tests

Introduction

- 14.1 The Wilcoxon Rank Sum Test
- 14.2 The Wilcoxon Signed Rank Test
- 14.3 The Kruskal-Wallis Test

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. We gave some practical guidelines for taking advantage of the robustness of these methods in Chapter 7.

robustness

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

outliers

1. If there are extreme **outliers** in a small data set, any inference method may be suspect. An outlier is an observation that may not come from the same population as the others. To decide what to do, you must find the cause of the outlier. 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,” it is risky to draw any conclusion from just a few observations. This is the advice we gave to the child development researcher in Example 2.19 (page 163).

transforming data

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. We used the logarithm transformation in Example 7.10 (page 519) to make the right-skewed distribution of carbon monoxide in vehicle exhausts more nearly normal.

3. In some settings, **other standard distributions** replace the normal distributions as models for the overall pattern in the population. We mentioned in Section 5.2 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 nonnormal models.

4. Finally, there are inference procedures that do not assume any specific form for the distribution of the population. These are called **nonparametric methods**. They are the subject of this chapter.

nonparametric methods

The word *nonparametric* contrasts these methods with statistical methods that are based on models of a specific form and use data to estimate the parameters in these models. For example, simple linear regression (Sections 2.3 and 10.1) uses a straight-line model. The parameters in this model are the

Setting	Normal test	Rank test
One sample	One-sample t test Section 7.1	Wilcoxon signed rank test Section 14.2
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 14.1
Several independent samples	One-way ANOVA F test Chapter 12	Kruskal-Wallis test Section 14.3

FIGURE 14.1 Comparison of tests based on normal distributions with nonparametric tests for similar settings.

slope and intercept of the line, and we can use the least-squares method to estimate these parameters from data. Scatterplot smoothers (Section 2.1), in contrast, do not assume any specific form for the relationship. When we use a smoother, we are doing **nonparametric regression**.

nonparametric
regression

This chapter concerns one type of nonparametric procedure, tests that can replace the t tests and one-way analysis of variance when the normality assumptions for those tests are not met. There are two big ideas that can serve as the basis for nonparametric tests. One is to use *counts*. This is the basis of the **sign test** for matched pairs, discussed on pages 519–522. The other is to use *ranks*. This chapter discusses rank tests.

sign test

Figure 14.1 presents an outline of the standard tests (based on normal distributions) and the rank tests that compete with them. All of these tests 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 14.1 all test hypotheses about population means. When distributions are strongly skewed, the mean may not be the preferred measure of center. We will see that rank tests do not test hypotheses about means.

We devote a section of this chapter to each of the rank procedures. Section 14.1, 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 14.2 and 14.3 more briefly describe other rank tests.

14.1 The Wilcoxon Rank Sum Test

Two-sample problems (see page 537) are among the most common in statistics. The most useful nonparametric significance test compares two distributions. Here is an example of this setting.

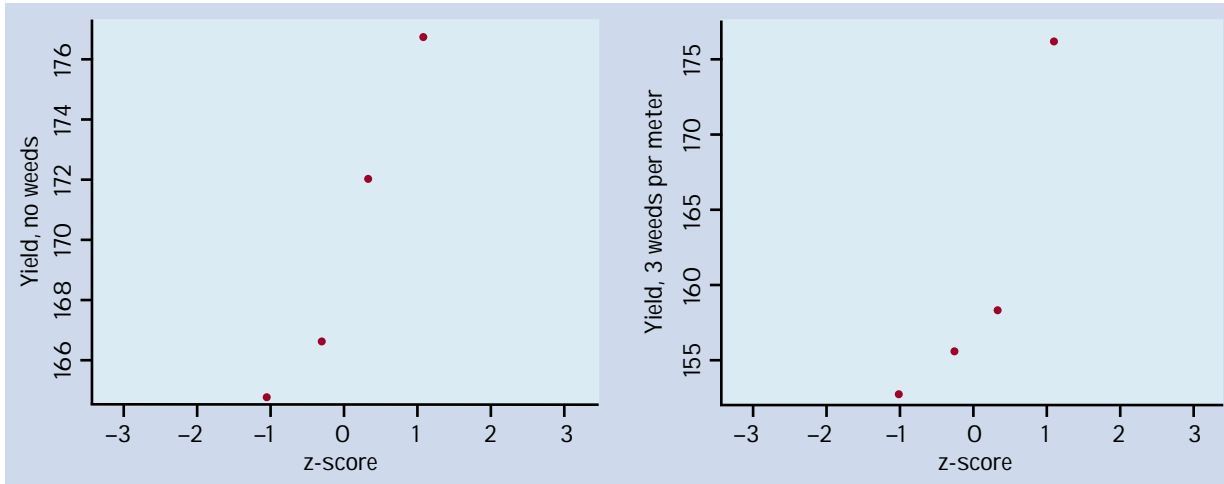


FIGURE 14.2 Normal quantile plots of corn yields from plots with no weeds (left) and with 3 weeds per meter of row (right).

EXAMPLE 14.1

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 8 small plots of ground, then weeded the corn rows by hand to allow no weeds in 4 randomly selected plots and exactly 3 lamb's-quarter plants per meter of row in the other 4 plots. Here are the yields of corn (bushels per acre) in each of the plots:¹

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

Normal quantile plots (Figure 14.2) suggest that the data 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 may prefer to use a test that does not require normality.

The rank transformation

We first rank all 8 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**

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	165.0	166.7	172.2	176.4	176.9
Rank	1	2	3	4	5	6	7	8

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.

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 either 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.*

In the corn yield study of Example 14.1, 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.

EXAMPLE 14.2

In Example 14.1, $n_1 = 4$, $n_2 = 4$, and there are $N = 8$ observations in all. The sum of ranks for the weed-free plots has mean

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

and standard deviation

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

Although the observed rank sum $W = 23$ is higher than the mean, it is only about 1.4 standard deviations high. 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.

*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 the American Cyanimid company.

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 packages give only approximate P -values. You must learn what your software offers.

EXAMPLE 14.3

Figure 14.3 shows the output from a software package 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 $P = 0.10$ 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.

It is worth noting that the two-sample t test gives essentially the same result as the Wilcoxon test in Example 14.3 ($t = 1.554$, $P = 0.0937$). It is in fact somewhat unusual to find a strong disagreement between the conclusions reached by these two 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 :

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

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.

continuity correction

```
Exact Wilcoxon rank-sum test
data: 0weeds and 3weeds
rank-sum statistic W = 23, n = 4, m = 4, p-value = 0.100
alternative hypothesis: true mu is greater than 0
```

FIGURE 14.3 Output from the S-Plus statistical software for the data in Example 14.1. This program uses the exact distribution for W when the samples are small and there are no ties (that is, when all observations have different values).

EXAMPLE 14.4

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 (see page 386) 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 $P = 0.10$.

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 of course safer for small samples. As Example 14.4 illustrates, however, the normal approximation with the continuity correction is often adequate.

EXAMPLE 14.5
Mann-Whitney test

Figure 14.4 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 alternate form of the Wilcoxon rank sum test. SAS carries out both the exact and 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 14.4. This approximate P -value is close to the exact result $P = 0.1000$, given by SAS and in Figure 14.3.

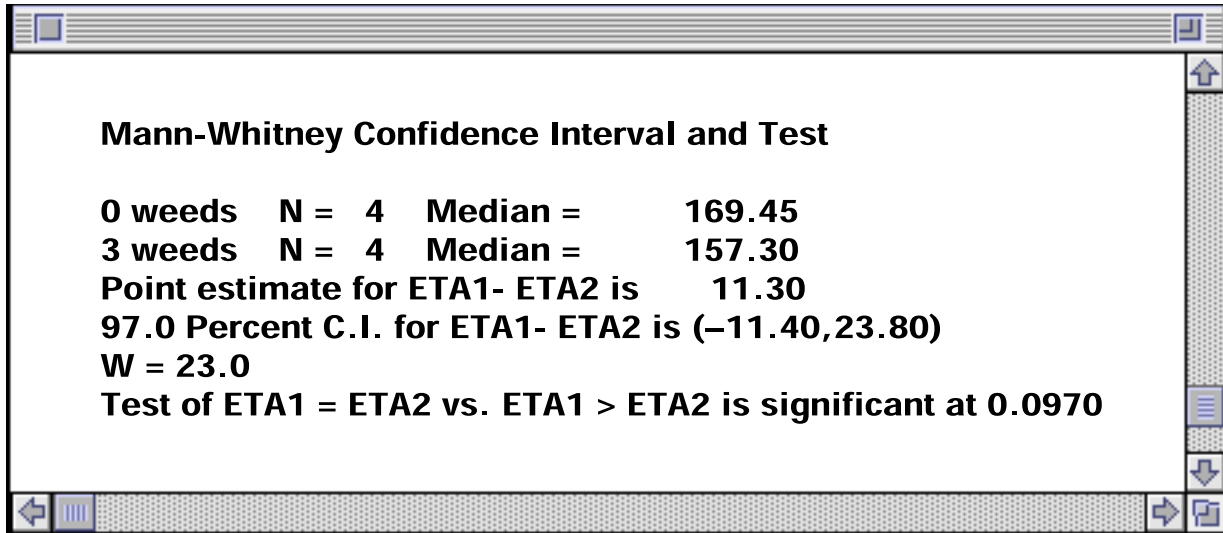
What hypotheses does Wilcoxon test?

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

$$\begin{aligned} H_0 : \mu_1 &= \mu_2 \\ H_a : \mu_1 &> \mu_2 \end{aligned}$$

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

$$\begin{aligned} H_0 : \text{median}_1 &= \text{median}_2 \\ H_a : \text{median}_1 &> \text{median}_2 \end{aligned}$$



(a)

Wilcoxon Scores (Rank Sums) for Variable YIELD
Classified by Variable WEEDS

WEEDS	N	Sum of Scores	Expected Under H0	Std Dev Under H0	Mean Score
0	4	23.0	18.0	3.46410162	5.75000000
3	4	13.0	18.0	3.46410162	3.25000000

Wilcoxon 2-Sample Test S = 23.0000

Exact P-Values
 (One-sided) Prob >= S = 0.1000
 (Two-sided) Prob >= |S - Mean| = 0.2000

Normal Approximation (with Continuity Correction of .5)
 Z = 1.29904 Prob > |Z| = 0.1939

(b)

FIGURE 14.4 Output from the Minitab and SAS statistical software for the data in Example 14.1. (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 provides a significance test for these hypotheses, but only if an additional assumption is met: both populations must have continuous distributions of *the same shape*. That is, the density curve for corn yields with 3 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 14.4(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. 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

$$\begin{aligned} H_0: & \text{two distributions are the same} \\ H_a: & \text{one has values that are systematically larger} \end{aligned}$$

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 3 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.²

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 6 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 only applies to data without ties. Moreover, the standard deviation σ_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 apply rank tests to data that have very many ties because the scale of measurement has only a few values. Here is an example.

EXAMPLE 14.6

Food sold at outdoor fairs and festivals may be less safe than food sold in restaurants because it is prepared in temporary locations and often by volunteer help. What do people who attend fairs think about the safety of the food served? One study asked this question of people at a number of fairs in the Midwest:

How often do you think people become sick because of food they consume prepared at outdoor fairs and festivals?

The possible responses were

- 1 = very rarely
- 2 = once in a while
- 3 = often
- 4 = more often than not
- 5 = always

In all, 303 people answered the question. Of these, 196 were women and 107 were men. Is there good evidence that men and women differ in their perceptions about food safety at fairs?³

We should first ask if the subjects in Example 14.6 are a random sample of people who attend fairs, at least in the Midwest. The researcher visited 11 different fairs. She stood near the entrance and stopped every 25th adult who passed. Because no personal choice was involved in choosing the subjects, we can reasonably treat the data as coming from a random sample. (As usual, there was some nonresponse, which could create bias.)

Here are the data, presented as a two-way table of counts:

	Response					Total
	1	2	3	4	5	
Female	13	108	50	23	2	196
Male	22	57	22	5	1	107
Total	35	165	72	28	3	303

Comparing row percentages shows that the women in the sample are more concerned about food safety than the men:

	Response					Total
	1	2	3	4	5	
Female	6.6%	55.1%	25.5%	11.7%	1.0%	100%
Male	20.6%	53.3%	20.6%	4.7%	1.0%	100%

Is the difference between the genders statistically significant?

We might apply the chi-square test (Chapter 9). It is highly significant ($\chi^2 = 16.120$, $df = 4$, $P = 0.0029$). Although the chi-square test answers our general question, it ignores the ordering of the responses and so does not use all of the available information. We would really like to know whether men or women are more concerned about food safety. This question depends on the ordering of responses from least concerned to most concerned. We can use the Wilcoxon test for the hypotheses

H_0 : men and women do not differ in their responses

H_a : one of the two genders gives systematically larger responses than the other

The alternative hypothesis is two-sided. Because the responses can take only five values, there are very many ties. All 35 people who chose “very rarely” are tied at 1, and all 165 who chose “once in a while” are tied at 2.

EXAMPLE 14.7

Figure 14.5 gives computer output for the Wilcoxon test. The rank sum for men (using average ranks for ties) is $W = 14,112.5$. The standardized value is $z = -3.259$ with two-sided P -value $P = 0.0011$. There is very strong evidence of a difference. Women are more concerned than men about the safety of food served at fairs.

Wilcoxon Scores (Rank Sums) for Variable SFAIR Classified by Variable GENDER					
GENDER	N	Sum of Scores	Expected Under H0	Std Dev Under H0	Mean Score
Female	196	31943.5000	29792.0	660.006327	162.977041
Male	107	14112.5000	16264.0	660.006327	131.892523
Average Scores Were Used for Ties					
Wilcoxon 2-Sample Test (Normal Approximation) (with Continuity Correction of .5)					
S = 14112.5 Z = -3.25906 Prob > Z = 0.0011					

FIGURE 14.5 Output from SAS for the food safety study in Example 14.6. The approximate two-sided P -value is 0.0011.

With more than 100 observations in each group and no outliers, we might use the two-sample t even though responses take only five values. In fact, the results for Example 14.6 are $t = 3.2719$ with $P = 0.0012$. The P -value for the two-sample t test is almost exactly equal to 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 “very rarely” and “once in a while” is the same as the difference between “once in a while” and “often.” 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 least to most concerned about safety, so the rank test makes sense. Some statisticians avoid using t procedures when there is not a fully meaningful scale of measurement.

Limitations of nonparametric tests

The examples we have given illustrate the potential usefulness of nonparametric tests. Nonetheless, rank tests are of secondary importance relative to inference procedures based on the normal distribution.

- Nonparametric inference 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 on to more advanced statistics.
- Normal tests compare means and are accompanied by simple confidence intervals for means or differences between means. When we use nonparametric tests to compare medians, we can also give confidence intervals, though they are rather awkward to calculate. However, the usefulness of nonparametric tests is clearest in settings when they do not simply compare medians—see the discussion of “What hypotheses does Wilcoxon test?” In these settings, there is no measure of the *size* of the observed effect that is closely related to the rank test of the *statistical significance* of the effect.
- The robustness of normal tests for means 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 14.6 because the response variable gives only the order of the responses.
- There are more modern and more effective ways to escape the assumption of normality, such as bootstrap methods (see page 445).

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 a 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 14.1 EXERCISES

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.

- 14.1** A study of early childhood education asked kindergarten students to tell a fairy tale 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 the children and assigned a score for certain uses of language. Here are the data (provided by Susan Stadler, Purdue University):

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

- Make normal quantile plots for the 5 responses in each group. Are any major deviations from normality apparent?
- 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?
- 14.2** Repeat the analysis of Exercise 14.1 for the scores when children retell a story they have heard and seen illustrated with pictures (Story 2).
- 14.3** Use the data in Exercise 14.1 for children telling Story 2 to carry out by hand the steps in the Wilcoxon rank sum test.
- (a) Arrange the 10 observations in order and assign ranks. There are no ties.
- (b) Find the rank sum W for the five 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?
- 14.4** The corn yield study of Example 14.1 also examined yields in four plots having 9 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 data above and some of the data from Example 14.1 to answer this question.
- (b) Compare the results from (a) with those from the two-sample t test for these data.
- (c) Now remove the low outlier 142.4 from the data with 9 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?
- 14.5** Example 7.17 (page 547) reports the results of a study of the effect of the pesticide DDT on nerve activity in rats. The data for the DDT group are

12.207	16.869	25.050	22.429	8.456	20.589
--------	--------	--------	--------	-------	--------

The control group data are

11.074	9.686	12.064	9.351	8.182	6.642
--------	-------	--------	-------	-------	-------

It is difficult to assess normality from such small samples, so we might use a nonparametric test to assess whether DDT affects nerve response.

- (a) State the hypotheses for the Wilcoxon test.
 - (b) Carry out the test. Report the rank sum W , its P -value, and your conclusion.
 - (c) The two-sample t test used in Example 7.17 found that $t = 2.9912$, $P = 0.0247$. Are your results different enough to change the conclusion of the study?
- 14.6** In Example 7.14, we compared the DRP scores of two groups of third graders who followed different reading curricula. The data appear in Table 7.2 (page 543).
- (a) Apply the Wilcoxon rank sum test to these data and compare your result with the $P = 0.0132$ obtained from the two-sample t test in Example 7.14.
 - (b) What are the null and alternative hypotheses for the t test? For the Wilcoxon test?
- 14.7** Table 7.3 (page 552) gives data on blood pressure before and after treatment for two groups of black males. One group took a calcium supplement, and the other group received a placebo. Example 7.20 compares the decrease in blood pressure in the two groups using pooled two-sample t procedures, and Exercise 7.69 applies the more general two-sample t procedures. The normal quantile plot for the calcium group (Figure 7.14, page 552) shows some departure from normality, though not enough to prevent use of t procedures.
- (a) Use the Wilcoxon rank sum test to analyze these data. Compare your findings with those of Example 7.20 (page 553) and Exercise 7.69 (page 563).
 - (b) What are the null and alternative hypotheses for each of the three tests we have applied to these data?
 - (c) What must we assume about the data to apply each of the three tests?
- 14.8** Exercise 7.51 (page 556) studies the effect of piano lessons on the spatial-temporal reasoning of preschool children. The data there concern 34 children who took piano lessons and a control group of 44 children. The data take only small whole-number values. Use the Wilcoxon rank sum test (there are many ties) to decide whether piano lessons improve spatial-temporal reasoning.
- 14.9** Example 14.6 describes a study of the attitudes of people attending outdoor fairs about the safety of the food served at such locations. The full data set is stored on the *IPS* CD as the file `eg14.06.dat`. It contains the responses of 303 people to several questions. The variables in this data set are (in order)

subject hfair sfair sfast srest gender

The variable “sfair” contains the responses described in the example concerning safety of food served at outdoor fairs. The variable “srest” contains responses to the same question asked about food served in restaurants. The variable “gender” contains 1 if the respondent is a woman, 2 if he is a man. We saw that women are more concerned than men about the safety of food served at fairs. Is this also true for restaurants?

- 14.10** The data file used in Example 14.6 and Exercise 14.9 contains 303 rows, one for each of the 303 respondents. Each row contains the responses of one person to several questions. We wonder if people are more concerned about the safety of food served at fairs than they are about the safety of food served at restaurants. Explain carefully why we *cannot* answer this question by applying the Wilcoxon rank sum test to the variables “sfair” and “srest.”
- 14.11** Shopping at secondhand stores is becoming more popular and has even attracted the attention of business schools. 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. (From William D. Darley, “Store-choice behavior for pre-owned merchandise,” *Journal of Business Research*, 27 (1993), pp. 17–31.)

Income code	Income	City 1	City 2
1	Under \$10,000	70	62
2	\$10,000 to \$19,999	52	63
3	\$20,000 to \$24,999	69	50
4	\$25,000 to \$34,999	22	19
5	\$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. The data file ex14.11.dat on the *IPS* CD contains data on the 459 shoppers in this study. The first variable is the city (City1 or City2) and the second is the income code as it appears in the table above (1 to 5). Is there good evidence that shoppers in one city have systematically higher incomes than in the other?

14.2 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. The **sign test** (page 519) is a nonparametric test

based on counts for matched pairs. We will now meet a rank test for this setting.

EXAMPLE 14.8

A study of early childhood education asked kindergarten students to tell a fairy tale 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:⁴

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 14.6) suggest a mild lack of normality. We would like to use a rank test.

Positive differences in Example 14.8 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

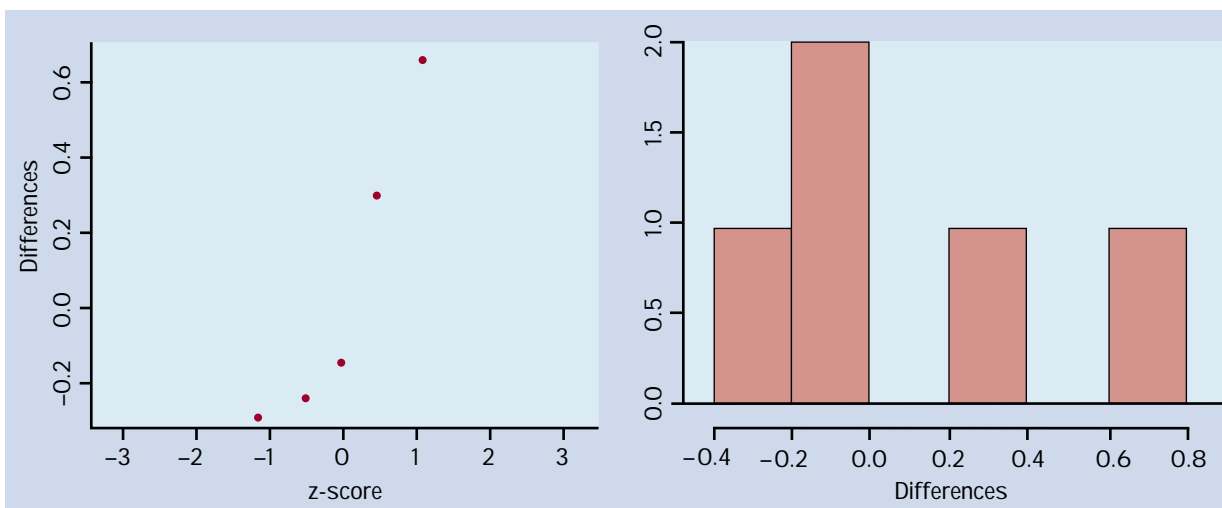


FIGURE 14.6 Normal quantile plot and histogram for the five differences in Example 14.8.

absolute value

absolute values of the differences, that is, their magnitudes without a sign. Here they are, with boldface indicating the positive values:

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 zero differences, discard them before ranking.

Absolute value	0.08	0.17	0.23	0.37	0.66
Rank	1	2	3	4	5

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 responses have a continuous distribution that 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.

EXAMPLE 14.9

In the storytelling study of Example 14.8, $n = 5$. 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 14.7 displays the output of two statistical programs. We see 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.

```

Exact Wilcoxon signed-rank test

data: Story2-Story1

signed-rank statistic V = 9, n = 5, p-value = 0.4062

alternative hypothesis: true mu is greater than 0

```

	Rank Totals	Cases	Mean Rank
Positive Ranks	9	2	4.500
Negative Ranks	6	3	2
Ties	•	0	•
Total	15	5	3

p = 0.4062

FIGURE 14.7 Output from (a) S-Plus and (b) Data Desk for the storytelling study in Example 14.9. These programs use the exact distribution of W^+ when the sample size is small and there are no ties.

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 14.10

For $n = 5$ observations, we saw in Example 14.9 that $\mu_{W^+} = 7.5$. The standard deviation of W^+ under the null hypothesis is

$$\sigma_{W^+} = \sqrt{\frac{n(n+1)(2n+1)}{24}} = \sqrt{\frac{(5)(6)(25)}{24}} = \sqrt{31.25} = 5.590$$

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:

$$P(W^+ \geq 8.5) = P\left(\frac{W^+ - 7.5}{5.590} \geq \frac{9 - 7.5}{5.590}\right) = P(Z \geq 0.27) = 0.394$$

Despite the small sample size, the normal approximation gives a result quite close to the exact value $P = 0.4062$.

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, we drop such pairs from our sample. As in the case of the Wilcoxon rank sum, ties complicate finding a P -value. There is no longer a usable exact distribution for the signed rank statistic W^+ , and the standard deviation σ_{W^+} must be adjusted for the ties before we can use the normal approximation. Software will do this. Here is an example.

EXAMPLE 14.11

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 : scores have the same distribution in Rounds 1 and 2

H_a : scores are systematically lower or higher in Round 2

A normal quantile plot of the differences (Figure 14.8) shows some irregularity and a low outlier. We will use the Wilcoxon signed rank test.

The absolute values of the differences, with boldface indicating those that are 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	3	3	4	5	5	5	5	5	6	16
Rank	1	2	3.5	3.5	5	8	8	8	8	8	11	12

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

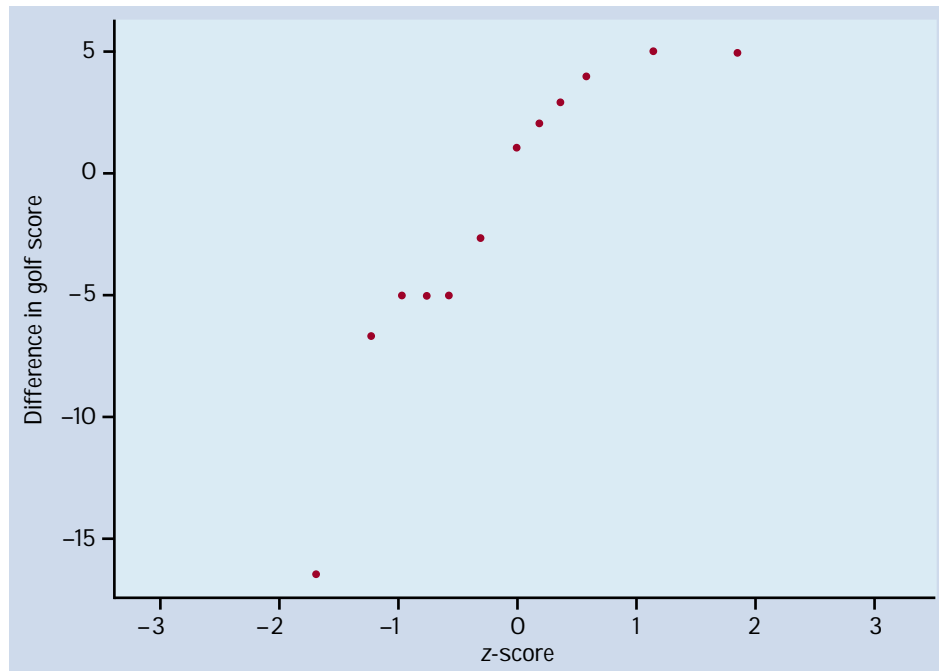


FIGURE 14.8 Normal quantile plot for Example 14.11 of the differences in scores for two rounds of a golf tournament.

EXAMPLE 14.12

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
Data Desk	$P = 0.366$
Minitab	$P = 0.388$
SAS	$P = 0.388$
S-PLUS	$P = 0.384$

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.

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 14.2 EXERCISES

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.

- 14.12** Table 7.1 (page 514) presents the scores on a test of understanding of spoken French for a group of high school French teachers before and after a summer language institute. The improvements in scores between the pretest and the posttest for the 20 teachers were

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

A normal quantile plot (Figure 7.7, page 515) shows granularity and a low outlier. We might wish to avoid the assumption of normality by using a rank test. Use the Wilcoxon signed rank procedure to reach a conclusion about the effect of the language institute. State hypotheses in words and report the statistic W^+ , its P -value, and your conclusion. (Note that there are many ties in the data.)

- 14.13** Exercise 7.35 (page 532) gives the scores on a test of comprehension of spoken Spanish for 20 teachers before and after they attended a summer language institute. We want to know whether the institute improves Spanish comprehension.
- State the null and alternative hypotheses.
 - Explain why the Wilcoxon rank sum test is *not* appropriate.
 - Give numerical measures that describe what the data show. Then use the Wilcoxon signed rank test to assess significance. What do you conclude?
- 14.14** Show the assignment of ranks and the calculation of the signed rank statistic W^+ for the data in Exercise 14.12. Remember that zeros are dropped from the data before ranking, so that n is the number of nonzero differences within pairs.
- 14.15** Example 14.6 describes a study of the attitudes of people attending outdoor fairs about the safety of the food served at such locations. The full data set is stored on the CD as the file eg14.06.dat. It contains the responses of 303 people to several questions. The variables in this data set are (in order):

subject hfair sfair sfast srest gender

The variable “sfair” contains responses to the safety question described in Example 14.6. The variable “srest” contains responses to the same question

asked about food served in restaurants. We suspect that restaurant food will appear safer than food served outdoors at a fair. Do the data give good evidence for this suspicion? (Give descriptive measures, a test statistic and its P -value, and your conclusion.)

- 14.16** The food safety survey data described in Example 14.6 and Exercise 14.15 also contain the responses of the 303 subjects to the same question asked about food served at fast-food restaurants. These responses are the values of the variable “sfast.” Is there a systematic difference between the level of concern about food safety at outdoor fairs and at fast-food restaurants?
- 14.17** Differences of electric potential occur naturally from point to point on a body’s skin. Is the natural electric field strength best for helping wounds to heal? If so, changing the field will slow healing. The research subjects are anesthetized newts. Make a razor cut in both hind limbs. Let one heal naturally (the control). Use an electrode to change the electric field in the other to half its normal value. After two hours, measure the healing rate. Here are healing rates (in micrometers per hour) for 14 newts. (Data provided by Drina Iglesia, Purdue University. The study results are reported in D. D. S. Iglesia, E. J. Cragoe, Jr., and J. W. Venable, “Electric field strength and epithelization in the newt (*Notophthalmus viridescens*),” *Journal of Experimental Zoology*, 274 (1996), pp. 56–62.)

Newt	Experimental limb	Control limb	Difference in healing
01	24	25	-1
02	23	13	10
03	47	44	3
04	42	45	-3
05	26	57	-31
06	46	42	4
07	38	50	-12
08	33	36	-3
09	28	35	-7
10	28	38	-10
11	21	43	-22
12	27	31	-4
13	25	26	-1
14	45	48	-3

The researchers want to know if changing the electric field reduces the healing rate for newts. State hypotheses, carry out a test, and give your conclusion. Be sure to include a description of what the data show in addition to the test results. (The researchers compared several field strengths and concluded that the natural strength is about right for fastest healing.)

- 14.18** Exercise 7.19 (page 527) presents these data on the vitamin C content of Haitian “bouillie” before and after cooking (in milligrams per 100 grams):

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

Is there a significant loss of vitamin C in cooking? Use a test that does not require normality.

- 14.19** Exercise 7.32 (page 531) gives data on the vitamin C content of 27 bags of wheat soy blend at the factory and five months later in Haiti. 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.
- 14.20** Exercise 7.33 (page 532) contains data from a student project that investigated whether right-handed people can turn a knob faster clockwise than they can counterclockwise. Describe what the data show, then state hypotheses and do a test that does not require normality. Report your conclusions carefully.

14.3 The Kruskal-Wallis Test

We have now considered alternatives to the matched pairs and two-sample t 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 14.13

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, 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:⁵

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 as follows:

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 14.9) show that outliers are present in the yields for 3 and 9 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 nonparametric test.

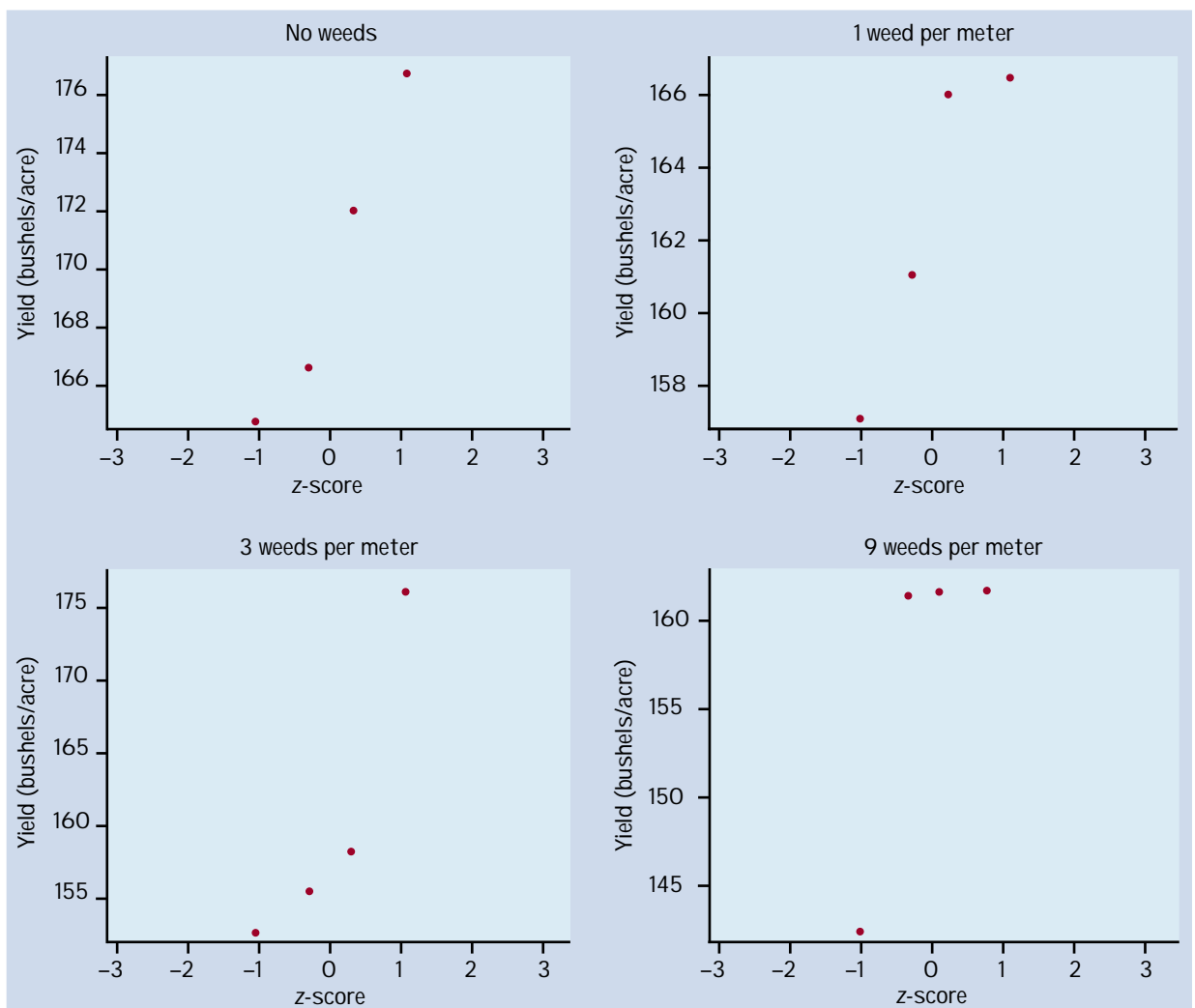


FIGURE 14.9 Normal quantile plots for the corn yields in the four treatment groups in Example 14.13.

Hypotheses and assumptions

The ANOVA F test concerns the means of the several populations represented by our samples. For Example 14.13, the ANOVA hypotheses are

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

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

For example, μ_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:

$$H_0: \text{yields have the same distribution in all groups}$$

$$H_a: \text{yields are systematically higher in some groups than in others}$$

“Systematically higher” has the precise meaning described in Section 14.1. 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, roughly speaking, 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 14.14

In Example 14.13, 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					Sum of ranks
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

Wilcoxon Scores (Rank Sums) for Variable YIELD
Classified by Variable WEEDS

WEEDS	N	Sum of Scores	Expected Under H0	Std Dev Under H0	Mean Score
0	4	52.5000000	34.0	8.24014563	13.1250000
1	4	33.5000000	34.0	8.24014563	8.3750000
3	4	25.0000000	34.0	8.24014563	6.2500000
9	4	25.0000000	34.0	8.24014563	6.2500000

Average Scores Were Used for Ties

Kruskal-Wallis Test (Chi-Square Approximation)
CHISQ = 5.5725 DF = 3 Prob > CHISQ = 0.1344

FIGURE 14.10 Output from SAS for the Kruskal-Wallis test applied to the data in Example 14.13. SAS uses the chi-square approximation to obtain a P -value.

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 G) 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 14.10 displays the output from the SAS statistical software, which gives the results $H = 5.5725$ and $P = 0.1344$. The software makes a small 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 14.14 is $P = 0.1299$. This result required several hours of computing time.⁶ 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.

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 14.3 EXERCISES

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

- 14.21** How do nematodes (microscopic worms) affect plant growth? A botanist prepares 16 identical planting pots and then introduces different numbers of nematodes into the pots. A tomato seedling is transplanted into each plot. Here are data on the increase in height of the seedlings (in centimeters) 16 days after planting. (Data provided by Matthew Moore.)

Nematodes	Seedling growth				
0	10.8	9.1	13.5	9.2	
1,000	11.1	11.1	8.2	11.3	
5,000	5.4	4.6	7.4	5.0	
10,000	5.8	5.3	3.2	7.5	

We applied ANOVA to these data in Exercise 12.10 (page 782). Because the samples are very small, it is difficult to assess normality.

- (a) What hypotheses does ANOVA test? What hypotheses does Kruskal-Wallis test?
- (b) Find the median growth in each group. Do nematodes appear to retard growth? Apply the Kruskal-Wallis test. What do you conclude?
- 14.22** The presence of harmful insects in farm fields is detected by erecting boards covered with a sticky material and then examining the insects trapped on the board. To investigate which colors are most attractive to cereal leaf beetles, researchers placed six boards of each of four colors in a field of oats in July. The table below gives data on the number of cereal leaf beetles

trapped. (Based on 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.)

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

We applied ANOVA to these data in Exercise 12.11 (page 782).

- (a) Make a normal quantile plot for each group. Are there indications of lack of normality?
 - (b) What hypotheses does ANOVA test? What hypotheses does Kruskal-Wallis test?
 - (c) 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?
- 14.23** Table 1.8 (page 40) presents data on the calorie and sodium content of selected brands of beef, meat, and poultry hot dogs. We will regard these brands as random samples from all brands available in food stores. We saw that the distribution of calories in meat hot dogs had two clusters and a low outlier. We might therefore prefer to use a nonparametric test. Give the five-number summaries for the three types of hot dog and then apply the Kruskal-Wallis test. Report your conclusions carefully.
- 14.24** Exercise 14.22 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.
- (a) What are I , the n_i , and N in this example?
 - (b) Arrange the counts in order and assign ranks. Be careful about ties. Find the sum of the ranks for each color.
 - (c) Calculate the Kruskal-Wallis statistic H . How many degrees of freedom should you use for the chi-square approximation of its null-hypothesis distribution? Use the chi-square table to give an approximate P -value.
- 14.25** Repeat the analysis of Exercise 14.23 for the sodium content of hot dogs, using the data in Table 1.8 (page 40).
- 14.26** Table 12.4 (page 781) gives data on the effect of four treatments on the spatial-temporal reasoning ability of preschool children. The treatments are piano lessons, singing lessons, computer instruction, and no lessons of any kind. The response variable is the change in a child's score on a test of spatial-temporal reasoning.
- (a) Give the five-number summary for each group. What do the data suggest about the effects of the treatments?

- (b) Make a normal quantile plot for each group. Which group deviates most from normality?
- (c) Do the treatments differ significantly in their ability to improve children's spatial-temporal reasoning?

14.27 Example 14.6 describes a study of the attitudes of people attending outdoor fairs about the safety of the food served at such locations. The full data set is stored on the CD as the file `eg14.06.dat`. It contains the responses of 303 people to several questions. The variables in this data set are (in order):

subject hfair sfair sfast srest gender

The variable “sfair” contains responses to the safety question described in Example 14.6. The variables “srest” and “sfast” contain responses to the same question asked about food served in restaurants and in fast-food chains. Explain carefully why we *cannot* use the Kruskal-Wallis test to see if there are systematic differences in perceptions of food safety in these three locations.

14.28 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. (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.) Here are the data, as a two-way table of counts:

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

The data for all 356 subjects are stored in the file `ex14.28.dat` on the CD. 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.

14.29 (Optional)
multiple comparisons

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.⁷ 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 14.23 you found a significant difference among the calorie contents of three types of hot dog. Now we will explore multiple comparisons.

- (a) Write down all of the pairwise comparisons we can make, for example, beef versus meat. There are three possible pairwise comparisons.
- (b) Carry out three Wilcoxon rank sum tests, one for each of the three pairs of hot dog types. 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?

14.30 (Optional)

Exercise 14.29 outlines how to use the Wilcoxon rank sum test several times for multiple comparisons with overall significance level 0.05 for all comparisons together. In Exercise 14.22 you found that the numbers of beetles attracted by boards of four colors differ significantly. At the overall 0.05 level, which pairs of colors differ significantly? (Hint: There are 6 possible pairwise comparisons among 4 colors.)

NOTES

1. Data provided by Sam Phillips, Purdue University. The data have been slightly modified to remove one tie to simplify Exercise 14.4.
2. 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.
3. Data from Huey Chern Boo, "Consumers' perceptions and concerns about safety and healthfulness of food served at fairs and festivals," M.S. thesis, Purdue University, 1997.
4. Data provided by Susan Stadler, Purdue University.
5. See Note 1.
6. Using SAS Version 6.12 on a 166 MHz Pentium personal computer (a quite fast machine at the time this was written) required about 3.5 hours to obtain the exact P -value in Example 14.14. No wonder very few software systems offer exact P -values for the Kruskal-Wallis statistic.
7. For more details on multiple comparisons (but not the simple procedure given here), see M. Hollander and D. A. Wolfe, *Nonparametric Statistical Methods*, Wiley, New York, 1973. This book is a useful reference on applied aspects of nonparametric inference in general.