

STATISTICS FOR DATA SCIENCE HYPOTHESIS and INFERENCE

Dr. Deepa NairDepartment of Science and Humanities



UNIT-4 HYPOTHESIS and INFERENCE

Session-8

Distribution Free Tests

Dr. Deepa Nair

Department of Science and Humanities

Distribution Free Tests.



- The samples are not required to come from any specific distribution.
- While distribution free tests do require assumptions for their validity, these assumptions are somewhat less restrictive than the assumptions needed for the *t* test.
- Distribution-free tests are sometimes called nonparametric tests.

Distribution Free Tests



• We discuss two distribution-free tests in this section. The first, called the Wilcoxon signed-rank test, is a test for a population mean, analogous to the one-sample *t* test.

 The second, called the Wilcoxon rank-sum test, or the Mann– Whitney test, is analogous to the two-sample t test.

Distribution Free Tests



- We illustrate this test with an example. The nickel content, in parts per thousand by weight, is measured for six welds. The results are 9.3, 0.9, 9.0, 21.7, 11.5, and 13.9. Let μ represent the mean nickel content for this type of weld. It is desired to test H_0 : $\mu \geq 12$ versus H_1 : $\mu < 12$.
- The Student's t test is not appropriate, because there are two outliers, 0.9 and 21.7, which indicate that the population is not normal. The Wilcoxon signed-rank test can be used in this situation. This test does not require the population to be normal.

Distribution Free Tests



The Wilcoxon Signed-Rank Test:

| X | X-12 | Rank |
|------|-------|------|
| 11.5 | -0.5 | -1 |
| 13.9 | 1.9 | 2 |
| 9.3 | -2.7 | -3 |
| 9.0 | -3.0 | -4 |
| 21.7 | 9.7 | 5 |
| 0.9 | -11.1 | -6 |

Distribution Free Tests



- Let $H_0: \mu \geq 12$, so a small value of S+ will provide evidence against H_0 .
- We observe S+=7. The P-value is the probability of observing a value of S+ that is less than or equal to 7 when H_0 is true.
- For sample size n=6, we find that the probability of observing a value of 4 or less is 0.1094.
- The probability of observing a value of 7 or less must be greater than this, so we conclude that P>0.1094, and thus do not reject H_0 .

Distribution Free Tests



The Wilcoxon Signed-Rank Test:

In the example discussed previously, the nickel content for six welds was measured to be 9.3, 0.9, 9.0, 21.7, 11.5, and 13.9. Use these data to test $H_0: \mu \leq 5$ versus $H_1: \mu > 5$.

Distribution Free Tests



The Wilcoxon Signed-Rank Test:

| X | X-5 | Rank |
|------|------|------|
| 11.5 | 6.5 | 4 |
| 13.9 | 8.9 | 5 |
| 9.3 | 4.3 | 3 |
| 9.0 | 4 | 1 |
| 21.7 | 16.7 | 6 |
| 0.9 | -4.1 | -2 |

Distribution Free Tests



- Let $H_0: \mu \leq 5$, so a large value of S+ will provide evidence against H_0 .
- We observe S+=19. The *P*-value is the probability of observing a value of S+ that is less than or equal to 7 when H_0 is true.
- For sample size n=6, the P- Value is the area in the right side of the null distribution corresponding to the values greater than or equal to 19.
- So we conclude that P value is 0.0469 < 0.05 and reject H_0 .

Distribution Free Tests



The Wilcoxon Signed-Rank Test:

In the example discussed previously, the nickel content for six welds was measured to be 9.3, 0.9, 9.0, 21.7, 11.5, and 13.9. Use these data to test H_0 : $\mu = 16$ versus H_1 : $\mu \neq 16$.

Distribution Free Tests



The Wilcoxon Signed-Rank Test:

| X | X-16 | Rank |
|------|-------|------|
| 11.5 | -4.5 | -2 |
| 13.9 | -2.1 | -1 |
| 9.3 | -6.7 | -4 |
| 9.0 | -7 | -5 |
| 21.7 | 5.7 | 3 |
| 0.9 | -15.1 | =6 |

Distribution Free Tests



- Let $H_0: \mu = 16$, this is a two tailed test.
- We observe S+=3.
- The *P*-value is the probability of observing a value of S + that is not equal to 3 when H_0 is true.
- So we conclude that P value is 0.1562 < 0.05 and reject H_0 .

Distribution Free Tests



Ties:

- Sometimes two or more of the quantities to be ranked have exactly the same value. Such quantities are said to be tied. The standard method for dealing with ties is to assign to each tied observation the average of the ranks they would have received if they had differed slightly.
- For example, the quantities 3, 4, 4, 5, 7
- Would receive the ranks 1, 2.5, 2.5, 4, 5
- The quantities 12, 15, 16, 16, 16, 20
- Would receive the ranks 1, 2, 4, 4, 4, 6.

Distribution Free Tests



Differences of Zero:

- If the mean under H_0 is μ_0 , and one of the observations is equal to μ_0 , then its difference is 0, which is neither positive nor negative.
- An observation that is equal to μ_0 cannot receive a signed rank. The appropriate procedure is to drop such observations from the sample altogether, and to consider the sample size to be reduced by the number of these observations.

Distribution Free Tests



Large-Sample Approximation:

- When the sample size *n* is large, the test statistic *S*+ is approximately normally distributed.
- A rule of thumb is that the normal approximation is good if n > 20.
- It can be shown by advanced methods that under H_0 , S+ has mean n(n+1)/4 and variance n(n+1)(2n+1)/24.

Distribution Free Tests



Large-Sample Approximation:

The z-score is

$$z = \frac{S_+ - n(n+1)/4}{\sqrt{(n+1)(2n+1)/24}}$$

Distribution Free Tests



Example 6.18

The article "Exact Evaluation of Batch-Ordering Inventory Policies in Two-Echelon Supply Chains with Periodic Review" (G. Chacon, Operations Research, 2001: 79–98) presents an evaluation of a reorder point policy, which is a rule for determining when to restock an inventory. Costs for 32 scenarios are estimated. Let μ represent the mean cost. Test $H_0: \mu \geq 70$ versus $H_1: \mu < 70$. The data, along with the differences and signed ranks, are presented in Table 6.1.

Distribution Free Tests



| X | X-70 | RANK | X | X-70 | RANK | Χ | X-70 | RANK |
|--------|--------|------|--------|--------|------|--------|--------|------|
| 79.26 | 9.26 | 1 | 30.27 | -39.70 | -12 | 11.48 | -58.52 | -23 |
| 80.79 | 10.79 | 2 | 22.39 | -47.61 | -13 | 11.28 | -58.72 | -24 |
| 82.07 | 12.07 | 3 | 118.39 | 48.39 | 14 | 10.08 | -59.92 | -25 |
| 82.14 | 12.14 | 4 | 118.46 | 48.46 | 15 | 7.28 | -62.72 | -26 |
| 57.19 | -12.81 | -5 | 20.32 | -49.68 | -16 | 6.87 | -63.13 | -27 |
| 55.86 | -14.14 | -6 | 16.69 | -53.31 | -17 | 6.23 | -63.77 | -28 |
| 42.08 | -27.92 | -7 | 16.50 | -53.50 | -18 | 4.57 | -65.43 | -29 |
| 41.78 | -28.22 | -8 | 15.95 | -54.05 | -19 | 4.09 | -65.91 | -30 |
| 100.01 | 30.01 | 9 | 15.16 | -54.84 | -20 | 140.09 | 70.09 | 31 |
| 100.36 | 30.36 | 10 | 14.22 | -55.78 | -21 | 140.77 | 70.77 | 32 |
| 30.46 | -39.54 | -11 | 11.64 | -58.36 | -22 | | | |

Distribution Free Tests



Solution:

The sample size is n=32, so the mean is n(n+1)/4=264 and the variance is n(n+1)(2n+1)/24=2860. The sum of the positive ranks is S+=121. We compute

$$z = \frac{121 - 264}{\sqrt{2860}} = -2.67$$

Since the null hypothesis is of the form H_0 : $\mu \geq \mu_0$, small values of S + provide evidence against H0. Thus the P-value is the area under the normal curve to the left of z=-2.67. This area, and thus the P-value, is 0.0038.

Distribution Free Tests



The Wilcoxon Rank-Sum Test:

- The Wilcoxon rank-sum test, also called the Mann– Whitney test, can be used to test the difference in population means in certain cases where the populations are not normal.
- Two assumptions are necessary.
- First the populations must be continuous.
- Second, their probability density functions must be identical in shape and size; the only possible difference between them being their location.

Distribution Free Tests



The Wilcoxon Rank-Sum Test:

- Let X_1, \ldots, X_m be a random sample from one population and let Y_1, \ldots, Y_n be a random sample from the other.
- We adopt the notational convention that when the sample sizes are unequal, the smaller sample will be denoted X_1, \ldots, X_m .
- Thus the sample sizes are m and n, with $m \leq n$.
- Denote the population means by μ_X and μ_Y , respectively.

Distribution Free Tests



The Wilcoxon Rank-Sum Test:

- The test is performed by ordering the m + n values obtained by combining the two samples, and assigning ranks 1, 2, ..., m + n to them.
- The test statistic, denoted by W, is the sum of the ranks corresponding to X_1, \ldots, X_m .

Distribution Free Tests



The Wilcoxon Rank-Sum Test:

- Since the populations are identical with the possible exception of location, it follows that if $\mu_X < \mu_Y$, the values in the X sample will tend to be smaller than those in the Y sample.
- So the rank sum W will tend to be smaller as well.
- By similar reasoning, if $\mu_X > \mu_Y$, W will tend to be larger.

Distribution Free Tests



The Wilcoxon Rank-Sum Test: Example:

• Resistances, in m, are measured for five wires of one type and six wires of another type. The results are as follows:

X: 36 28 29 20 38

Y: 34 41 35 47 49 46

• Use the Wilcoxon rank-sum test to test $H_0: \mu_X \ge \mu_Y \ versus H_1: \mu_X < \mu_Y$.

Distribution Free Tests

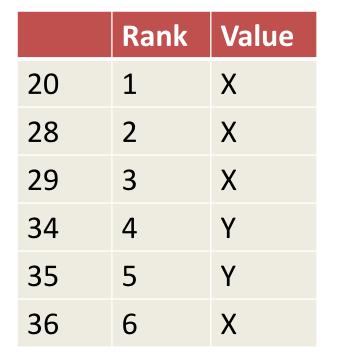


The Wilcoxon Rank-Sum Test: Solution:

We order the 11 values and assign the ranks.

Distribution Free Tests

The Wilcoxon Rank-Sum Test: Solution:



| | Rank | Value |
|----|------|-------|
| 38 | 7 | X |
| 41 | 8 | Υ |
| 46 | 9 | Υ |
| 47 | 10 | Υ |
| 49 | 11 | Υ |



Distribution Free Tests



The Wilcoxon Rank-Sum Test: Solution:

$$W = 1 + 2 + 3 + 6 + 7 = 19.$$

- To determine the P-value, we consult Table A.6 (in Appendix A).
- We note that small values of W provide evidence against H_0 : $\mu_X \geq \mu_Y$, so the P value

Distribution Free Tests

The Wilcoxon Rank-Sum Test: Solution:

• Is the area in the left-hand tail of the null distribution. Entering the table with $m=5\ and\ n=6$ we find that the area to the left of W=19 is 0.0260. This is the P-value.



Distribution Free Tests



Large-Sample Approximation:

- When both sample sizes m and n are greater than 8, it can be shown by advanced methods that the null distribution of the test statistic W is approximately normal with mean m(m+n+1)/2 and variance mn(m+n+1)/12.
- z score is $z = \frac{W m(m + n + 1)/2}{\sqrt{mn(m + n + 1)/12}}$

Distribution Free Tests



Example:

The article "Cost Analysis Between SABER and Design Bid Build Contracting Methods" (E. Henry and H. Brothers, Journal of Construction Engineering and Management, 2001:359–366) presents data on construction costs for 10 jobs bid by the traditional method (denoted X) and 19 jobs bid by an experimental system (denoted Y). The data, in units of dollars per square meter, and their ranks, are presented in Table 6.2. Test H_0 : $\mu_X \leq \mu_Y versus\ H_1: \mu_X > \mu_Y$.

Distribution Free Tests



| VALUE | RANK | SAMPLE | VALUE | RANK | SAMPLE |
|-------|------|--------|-------|------|--------|
| 57 | 1 | X | 613 | 16 | X |
| 95 | 2 | Υ | 622 | 17 | Y |
| 101 | 3 | Y | 708 | 18 | X |
| 118 | 4 | Υ | 726 | 19 | Y |
| 149 | 5 | Y | 843 | 20 | Y |
| 196 | 6 | Υ | 908 | 21 | Υ |
| 200 | 7 | Υ | 926 | 22 | X |
| 233 | 8 | Υ | 943 | 23 | Y |
| 243 | 9 | Y | 1048 | 24 | Υ |

Distribution Free Tests



| VALUE | RANK | SAMPLE | VALUE | RANK | SAMPLE |
|-------|------|--------|-------|------|--------|
| 341 | 10 | X | 1165 | 25 | X |
| 419 | 11 | Υ | 1293 | 26 | X |
| 457 | 12 | X | 1593 | 27 | X |
| 584 | 13 | X | 1952 | 28 | X |
| 592 | 14 | Y | 2424 | 29 | Υ |
| 594 | 15 | Y | | | |

Distribution Free Tests



Solution:

The sum of the X ranks is W=1+12+13+16+18+22+25+26+27+28=188. The sample sizes are m=10 and n=19. We use the normal approximation and compute

$$z = \frac{\frac{188 - 10(10 + 19 + 1)}{2}}{\frac{\sqrt{10(19)(10 + 19 + 1)}}{12}}$$
= 1.74

Large values of W provide evidence against the null hypothesis. Therefore the P-value is the area under the normal curve to the right of z=1.74. From the z table we find that the P-value is 0.0409

Distribution Free Tests



Distribution-Free Methods Are Not Assumption-Free

We have pointed out that the distribution-free methods presented here require certain assumptions for their validity. Unfortunately, this is sometimes forgotten in practice. It is tempting to turn automatically to a distribution-free procedure in any situation in which the Student's t test does not appear to be justified, and to assume that the results will always be valid. This is not the case. The necessary assumptions of symmetry for the signed-rank test and of identical shapes and spreads for the rank-sum test are actually rather restrictive. While these tests perform reasonably well under moderate violations of these assumptions, they are not universally applicable.



Dr. Deepa Nair

Department of Science and Humanities

deepanair@pes.edu