

# **STATISTICS FOR DATA SCIENCE HYPOTHESIS and INFERENCE**

**Dr. Deepa Nair**Department 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**



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



#### **Distribution Free Tests**



## The Wilcoxon Signed-Rank Test:

## **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  $\mu$  represent the mean nickel content for this type of weld.

#### **Distribution Free Tests**



- 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.

#### **Distribution Free Tests**



- To compute the rank-sum statistic, we begin by subtracting 12 from each sample observation to obtain differences. The difference closest to 0, ignoring sign, is assigned a rank of 1.
- The difference next closest to 0, again ignoring sign, is assigned a rank of 2, and so on.
- Finally, the ranks corresponding to negative differences are given negative signs. The following table shows the results.

## **Distribution Free Tests**



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**



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**



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  $\mu_0$ , and one of the observations is equal to  $\mu_0$ , then its difference is 0, which is neither positive nor negative.
- An observation that is equal to  $\mu_0$  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**



- 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**



- 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  $\mu_X$  and  $\mu_Y$  , respectively.

#### **Distribution Free Tests**



- 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**



- 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**

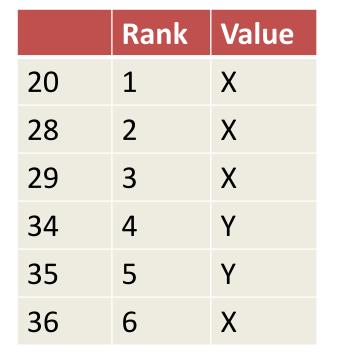


The Wilcoxon Rank-Sum Test: Solution:

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**

## PES UNIVERSITY ONLINE

## 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.
- $z-score\ is$   $z = \frac{W m(m+n+1)/2}{\sqrt{mn(m+n+1)/12}}$



## Dr. Deepa Nair

Department of Science and Humanities

deepanair@pes.edu