

Learning the Wilcoxon Signed-Rank Test: A Comprehensive Guide

Authored by
Mohammed loot

November 8, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *Learning the Wilcoxon Signed-Rank Test: A Comprehensive Guide*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=13841>

Introduction to the Wilcoxon Signed Rank Test

The [Wilcoxon Signed Rank Test](#) (WSRT) is a foundational procedure within [non-parametric](#) statistics. It serves as the definitive alternative to the standard [paired t-test](#), specifically when researchers encounter data that fail to satisfy the strict distributional assumptions of parametric methods. This test is meticulously engineered for analyzing dependent samples or repeated measures data, aiming to assess whether a statistically significant shift exists in the central tendency (typically the median) between two related measurements taken from the same subjects.

Crucially, unlike the paired t-test, the WSRT does not require the differences between the paired observations to conform to a normal [distribution](#). This resilience makes the Wilcoxon test exceptionally valuable and **robust** when dealing with data that is severely skewed, originates from small samples, or is measured on an ordinal scale. Its methodology provides a reliable path for statistical inference when traditional techniques would compromise the validity of the conclusions.

The underlying principle of the WSRT involves a thoughtful examination of the differences observed between each pair. It accounts for both the **magnitude** and the **direction** (sign) of these differences. By ranking the absolute differences, the procedure naturally assigns greater weight to observations that exhibit larger discrepancies, offering a comprehensive and sensitive assessment of any systematic shift in score distribution across the two measurement points. Mastering the Wilcoxon Signed Rank Test is therefore essential for any researcher whose data frequently falls outside the stringent requirements of parametric analysis.

Determining When to Apply the Wilcoxon Signed Rank Test

The primary factor that mandates the use of the [Wilcoxon Signed Rank Test](#) is the inability to satisfy the core assumptions underpinning the [paired t-test](#). Researchers who initially plan to utilize the paired t-test must first rigorously verify that the distribution of the differences calculated between the pairs is reasonably normal. When this assumption is severely violated--meaning the distribution of differences is distinctly non-normal--the reliability of the t-test results, particularly when working with smaller sample sizes, is seriously compromised, necessitating the shift to this non-parametric solution.

A fundamental diagnostic tool for assessing normality is the visual inspection of a **histogram** plotted from the calculated differences. If this graphical representation deviates dramatically from the classic symmetrical, "bell-shaped" curve characteristic of a normal [distribution](#), the use of the Wilcoxon test is statistically justified. However, analysts must recognize that the paired t-test is often quite **robust** against minor, subtle departures from perfect normality. Consequently, the observed deviation must be notably severe--perhaps demonstrating extreme skewness, pronounced heavy tails, or clear bimodal characteristics--to warrant abandoning the inherent

statistical power of the paired t-test.

While the WSRT sacrifices a marginal amount of statistical power compared to the paired t-test when the assumption of normality holds true, it offers substantially greater reliability and validity when that assumption is clearly absent. Therefore, the decision to pivot to the Wilcoxon test is not taken lightly but is instead a critical decision to maintain statistical integrity in the face of problematic data distributions.

Practical Application: A Training Program Example

To clearly illustrate the execution of the [Wilcoxon Signed Rank Test](#), we will analyze a common scenario involving athletic performance improvement. Imagine a basketball coach who introduces a new, specialized training regimen and wishes to statistically confirm whether this program genuinely leads to an increase in the number of successful free throws made by his players. To collect the required data, the coach measures the performance of 15 athletes, recording the total number of free throws made (out of 20 attempts) immediately before the training program starts and again following its completion.

This experimental design inherently produces **paired data**, as each player's "pre" score is fundamentally linked to their corresponding "post" score. The coach initially planned to use a [paired t-test](#) to evaluate the mean difference in performance. However, upon plotting the distribution of the score differences (Post score minus Pre score), it was definitively determined that this [distribution](#) was significantly **non-normal**. Given this distributional irregularity, relying on the parametric t-test would violate established statistical principles. Consequently, the coach correctly chose the more appropriate Wilcoxon Signed Rank Test to draw a reliable and trustworthy conclusion regarding the program's effectiveness.

The subsequent detailed steps outline the precise manual calculations required to execute this robust analysis using the raw performance data collected from the 15 athletes.

The following table summarizes the observed performance metrics--the number of free throws made (out of 20 attempts)--for each of the 15 players, both before and after the dedicated training program:

Player	Before	After
Player #1	14	15
Player #2	17	17
Player #3	12	15
Player #4	15	15
Player #5	15	17
Player #6	9	14
Player #7	12	9
Player #8	13	14
Player #9	13	11
Player #10	15	16
Player #11	19	18
Player #12	17	20
Player #13	14	20
Player #14	14	10
Player #15	16	17

Step-by-Step Procedure for the Wilcoxon Signed Rank Test

The successful execution of the Wilcoxon Signed Rank Test demands a systematic, four-stage process involving precise data transformation, meticulous ranking, and final summation of the weighted ranks.

Step 1: State the Null and Alternative Hypotheses

Establishing the hypotheses provides the essential statistical framework for inference. The focus of the WSRT is centered on the median difference, denoted as η_D .

H_0 : The median difference (η_D) between the two groups is zero. (i.e., The training program has no measurable effect on free throw performance.)

H_A : The median difference (η_D) is negative. (The post-training scores are higher than the pre-training scores, indicating an improvement caused by the program. This assumes the difference is calculated as Pre minus Post.)

Step 2: Find the Difference and Absolute Difference for Each Pair

For every player, the score difference is calculated by subtracting the Post-Training Score from the Pre-Training Score. This difference value must retain its original sign (positive or negative). Immediately after this, the **absolute difference** is determined by taking the magnitude of the difference, intentionally disregarding the sign. This absolute value forms the fundamental basis for

the subsequent ranking process, as the test is initially focused solely on the extent of the change, regardless of direction.

Player	Before	After	Difference	Abs. Difference
Player #1	14	15	-1	1
Player #2	17	17	0	0
Player #3	12	15	-3	3
Player #4	15	15	0	0
Player #5	15	17	-2	2
Player #6	9	14	-5	5
Player #7	12	9	3	3
Player #8	13	14	-1	1
Player #9	13	11	2	2
Player #10	15	16	-1	1
Player #11	19	18	1	1
Player #12	17	20	-3	3
Player #13	14	20	-6	6
Player #14	14	10	4	4
Player #15	16	17	-1	1

Step 3: Rank the Absolute Differences

This step is the core non-parametric operation. It involves ordering all pairs based on their absolute differences, assigning a rank from 1 (smallest absolute difference) to n (largest). Crucially, any pair exhibiting an absolute difference of exactly "0" must be **ignored** entirely and excluded from the effective sample size (n). If multiple pairs share the exact same absolute difference (a common occurrence known as a tie), the standard procedure dictates assigning the **mean rank** of the tied positions to each of those tied pairs. This technique ensures a consistent and statistically unbiased assignment of weight across all observations.

Player	Before	After	Difference	Abs. Difference	Rank
Player #2	17	17	0	0	-
Player #4	15	15	0	0	-
Player #1	14	15	-1	1	3
Player #8	13	14	-1	1	3
Player #10	15	16	-1	1	3
Player #11	19	18	1	1	3
Player #15	16	17	-1	1	3
Player #5	15	17	-2	2	6.5
Player #9	13	11	2	2	6.5
Player #3	12	15	-3	3	9
Player #7	12	9	3	3	9
Player #12	17	20	-3	3	9
Player #14	14	10	4	4	11
Player #6	9	14	-5	5	12
Player #13	14	20	-6	6	13

Step 4: Sum the Signed Ranks

Once the ranks are assigned, the original sign of the difference (from Step 2) is carefully reapplied to the corresponding rank. This action segregates the ranks into two critical components: **Positive Ranks** (where the Post score was lower or the Pre score was higher, indicating a negative change) and **Negative Ranks** (where the Post score was higher, indicating a positive change or improvement). Finally, the sums of these two groups are calculated independently, yielding the sum of positive ranks ($\sum R^+$) and the sum of negative ranks ($\sum R^-$).

Player	Before	After	Difference	Abs. Difference	Rank	Negative Ranks	Positive Ranks
Player #2	17	17	0	0	-		
Player #4	15	15	0	0	-		
Player #1	14	15	-1	1	3	-3	
Player #8	13	14	-1	1	3	-3	
Player #10	15	16	-1	1	3	-3	
Player #11	19	18	1	1	3		3
Player #15	16	17	-1	1	3	-3	
Player #5	15	17	-2	2	6.5	-6.5	
Player #9	13	11	2	2	6.5		6.5
Player #3	12	15	-3	3	9	-9	
Player #7	12	9	3	3	9		9
Player #12	17	20	-3	3	9	-9	
Player #14	14	10	4	4	11		11
Player #6	9	14	-5	5	12	-12	
Player #13	14	20	-6	6	13	-13	
Sum						-61.5	29.5

Calculating the Test Statistic and Critical Value

The **test statistic**, conventionally denoted as W , is the definitive output of the [Wilcoxon Signed Rank Test](#). W is mathematically defined as the smaller of the absolute values of the two rank sums ($\sum R^+$ and $\sum R^-$). In the context of our free throw performance data, the calculated rank sums are directly compared, and the minimum value is identified. In this specific worked example, the smaller sum is 29.5. Consequently, our calculated test statistic is $W = 29.5$. This statistic quantifies how far the observed data deviates from the expectations set forth by the [null hypothesis](#).

To establish statistical significance, the calculated W statistic must be rigorously compared against the [critical value](#). This vital threshold is located by consulting a standard Wilcoxon Signed Rank distribution table, which necessitates two key inputs: the effective sample size (n) and the predetermined [alpha level](#) (α). The effective sample size, n , is calculated as the total number of paired observations minus any pairs that were intentionally ignored because they yielded a difference of zero. Starting with 15 players and excluding two pairs with zero differences, our effective sample size is $n = 13$.

Assuming a conventional [alpha level](#) of $\alpha = 0.05$ (for a two-tailed test, or a directional test depending on the table used), we consult the critical value table using $n=13$. The corresponding [critical value](#) for this specific scenario is found to be 17 . This value functions as the strict rejection threshold: if the calculated test statistic W is less than or equal to this value, the observed result is deemed statistically rare enough to warrant rejection of the [null hypothesis](#).

n	Alpha value				
	0.005	0.01	0.025	0.05	0.10
5	-	-	-	-	0
6	-	-	-	0	2
7	-	-	0	2	3
8	-	0	2	3	5
9	0	1	3	5	8
10	1	3	5	8	10
11	3	5	8	10	13
12	5	7	10	13	17
13	7	9	13	17	21
14	9	12	17	21	25
15	12	15	20	25	30
16	15	19	25	29	35
17	19	23	29	34	41
18	23	27	34	40	47
19	27	32	39	46	53
20	32	37	45	52	60
21	37	42	51	58	67
22	42	48	57	65	75
23	48	54	64	73	83
24	54	61	72	81	91
25	60	68	79	89	100
26	67	75	87	98	110
27	74	83	96	107	119
28	82	91	105	116	130
29	90	100	114	126	140
30	98	109	124	137	151

Interpreting the Results and Drawing Conclusions

The final and most crucial phase of the Wilcoxon Signed Rank Test involves interpreting the relationship between the calculated test statistic (T) and the established critical value. The decision rule is unambiguous: if T is *less than or equal to* the critical value, the **null hypothesis** must be rejected. This outcome conclusively suggests that the observed difference is statistically significant. Conversely, if T is greater than the critical value, we formally fail to reject the null hypothesis, concluding that insufficient evidence exists to support the claim of a statistically significant effect.

Applying this rigorous rule to our basketball training example, we compare our calculated test statistic of $T = 29.5$ against a critical value of 17 (based on $n=13$ and $\alpha=0.05$). Since 29.5 is substantially larger than 17 , we are compelled to **fail to reject the null hypothesis**. This critical statistical finding indicates that, despite any numerical improvements observed in the raw free throw performance data, we do not possess sufficient evidence at the 5% **alpha level** to

confidently state that the training program leads to a statistically significant increase in performance. The results suggest the training program, as implemented and measured, was not statistically proven to be effective in shifting the median performance score.

While the manual execution of the Wilcoxon Signed Rank Test offers profound insight into its mechanics, its inherent complexity--particularly concerning the precise handling of ties and the calculation of ranks--makes it highly susceptible to errors, especially when dealing with larger datasets. Utilizing modern statistical software is nearly always the preferred methodology for ensuring accuracy, efficiency, and obtaining the precise p-value, which offers a more granular conclusion than relying solely on the fixed critical value tables.

Note: Use the

R Wilcoxon Signed Rank Test function

if you wish to perform the test using statistical software or a calculator instead of by hand. Automated statistical packages streamline the calculation of the test statistic and the p-value, offering a more precise and reliable conclusion than relying solely on critical value tables.