

How to Perform a Wilcoxon Signed-Rank Test to Compare Related Samples

Authored by
stats writer

January 22, 2026

RECOMMENDED CITATION

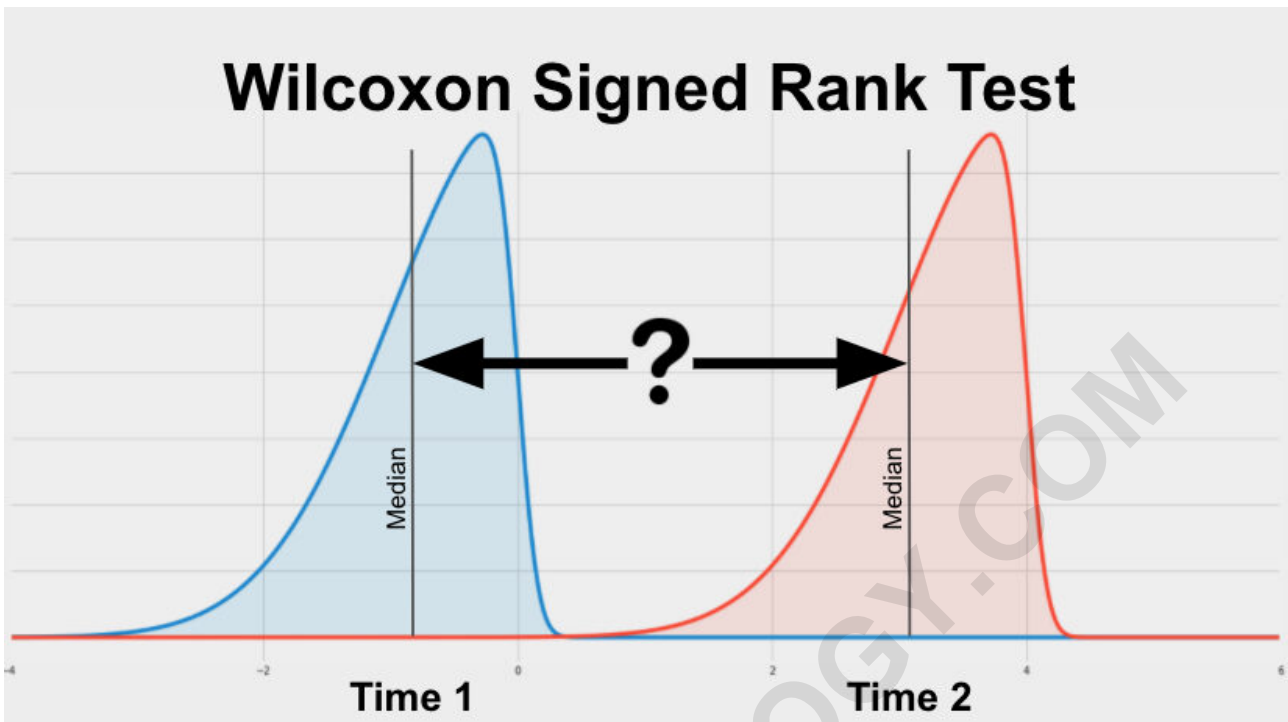
stats writer (2026). *How to Perform a Wilcoxon Signed-Rank Test to Compare Related Samples*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=126986>

The Wilcoxon Signed-Rank Test is a cornerstone non-parametric statistical test designed specifically for analyzing related samples. Its primary function is to assess whether two dependent measurements originate from the same population distribution, or if a significant shift has occurred between them. This powerful technique is indispensable when data fails to satisfy the stringent prerequisites of parametric alternatives, such as the t-test, particularly the assumption of normality. The methodology centers on calculating the differences between paired observations, ranking the absolute values of these differences, and then summing the ranks to derive a test statistic. This statistic provides the basis for determining if the observed difference is statistically significant, offering crucial insights into the effectiveness of interventions or the relationship between two variables in research settings.

What is the Wilcoxon Signed-Rank Test?

The **Wilcoxon Signed-Rank Test** (WSRT) serves as a robust statistical tool used primarily to compare two sets of measurements taken from the same group or from matched individuals. This test is specifically designed to determine if there is a statistically significant difference between these two related measurements on a continuous variable of interest. It is essential for researchers studying 'before-and-after' scenarios, repeated measures designs, or comparisons involving naturally paired data points.

Unlike its parametric counterpart, the paired samples t-test, the WSRT operates effectively even when the underlying data distributions are skewed or non-normal. To properly utilize this test, the variable under investigation must be continuous--meaning it can take on any value within a given range--and the sample from which the data is drawn must be obtained via simple random sample selection. Adhering to these structural assumptions ensures the validity and generalizability of the statistical inferences drawn from the analysis.



The Wilcoxon Signed-Rank Test is also widely recognized under several alternative names, reflecting its specific application in various fields. These include the Matched Pairs Wilcoxon Test, the Wilcoxon T-Test, the Wilcoxon Sign Test, and the highly descriptive Wilcoxon Matched Pairs Signed-Rank Test. Regardless of the nomenclature used, its fundamental application remains the comparison of two dependent measures.

Critical Assumptions for Applying the Wilcoxon Signed-Rank Test

Every statistical method, whether parametric or non-parametric, is built upon a foundation of underlying assumptions. These assumptions dictate the necessary characteristics of your dataset for the statistical results to be considered accurate, reliable, and interpretable. Failing to meet these foundational criteria can lead to invalid conclusions or an inaccurate measure of effect. The Wilcoxon Signed-Rank Test, being a non-parametric method, has fewer and less restrictive assumptions than, for example, the paired t-test, making it a highly adaptable choice for many real-world datasets.

Understanding and verifying these prerequisites is a crucial step in the analytical process. When planning a study or preparing data for analysis, researchers must ensure their variables and sampling methodology align with these requirements to maximize the integrity of the findings. The core assumptions specific to the WSRT are summarized as follows:

The variable of interest must be ****Continuous**** or ordinal.

The distribution of the differences between the paired measurements must be **Symmetric** (though the individual samples can be **Skewed**).

The data must be collected using a **Random Sample** procedure.

To provide a comprehensive understanding of when and how to apply the WSRT correctly, let us explore each of these critical assumptions in greater detail, focusing on their practical implications for research data.

Data Level of Measurement: Continuous Variables

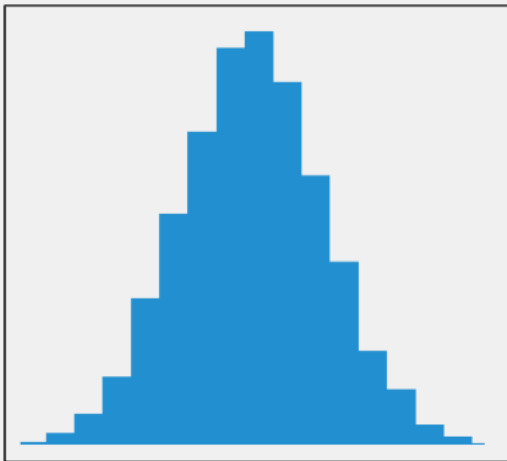
The dependent variable--the measure you are comparing between the two paired observations--must be measured on a continuous scale. A continuous variable is defined as one that can theoretically take on any value within a range, offering high precision in measurement. This high level of measurement allows for meaningful calculations of differences and rankings, which are central to the WSRT methodology.

Classic examples of suitable continuous variables include quantifiable measures such as age measured precisely in years, weight, height, standardized test scores, or psychological survey scores derived from aggregated items. These examples highlight data where fine distinctions between values are possible and meaningful. If your variable is not continuous but rather consists of frequencies or proportions, the WSRT is not the correct choice. For instance, if you are comparing two proportions (e.g., the percentage of voters favoring Candidate A before and after a debate), you would likely need to employ the McNemar Test instead.

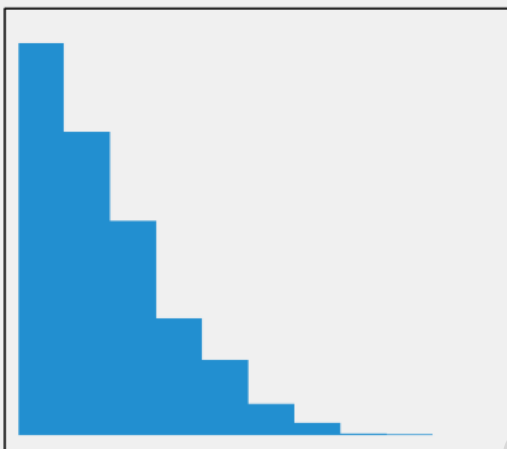
Distribution Shape Flexibility: Accommodating Skewness

One of the primary benefits of using a non-parametric test like the WSRT is its flexibility regarding the distributional shape of the data. The variable of interest does not need to conform to a bell-shaped curve, which is technically known as being normally distributed. This relaxation of the normality assumption makes the Wilcoxon test highly valuable when dealing with datasets that naturally exhibit pronounced skewness, heavy tails, or other deviations from normality.

However, it is important to note a subtle but crucial requirement regarding the distribution of the *differences* between the paired scores: while the individual sample distributions can be highly skewed, the distribution of the differences must be symmetrical around the median difference. If the data is substantially skewed, researchers are encouraged to use the Wilcoxon Signed-Rank Test, thereby avoiding the potential for inaccurate results that can arise if a Paired Samples T-Test is applied to non-normal data.



A normal distribution.
It is bell shaped with most of the data in the middle



A skewed distribution.
It is leaning left or right with most of the data on the edge

If, upon inspection, your variable of interest is confirmed to be approximately normally distributed, the appropriate choice for comparison between two paired groups would revert to the Paired Samples T-Test, as parametric tests often offer greater statistical power when their assumptions are met.

Sampling Method: Importance of Randomization

A fundamental requirement for drawing valid statistical inferences--that is, generalizing findings from a sample to a larger population--is that the data points must be derived from a simple random sample. This procedural assumption mandates that every individual in the target population has an equal chance of being selected for the study. When employing the WSRT in a paired design, the underlying group from which the measurements are taken must be randomly selected.

Consider a scenario where researchers wish to test the efficacy of a new diet plan. They must randomly select participants from the population of interest to form the study group. The key is ensuring that the selection mechanism itself is unbiased. If groups are determined based on convenience, self-selection, or any non-random criteria, the resulting analysis may suffer from

statistical bias, leading to incorrect or misleading conclusions that cannot be reliably extended beyond the immediate sample. While statistical tests can still be computed without random sampling, the conclusions are severely limited in scope.

If you cannot achieve a random sample, the generalizability of your results is inherently limited, restricting your ability to make strong population inferences. If your study design involves comparing two entirely independent samples (e.g., Group A receiving Treatment X vs. Group B receiving Placebo), rather than paired samples, you should utilize the Mann-Whitney U Test instead.

Interpreting Results: The Condition of Similar Shape

While the WSRT can accommodate skewed data, the interpretation of a significant result is dependent on the morphological similarity of the two measured distributions. To confidently assert that a significant Wilcoxon result indicates a difference in central tendency--specifically, a difference in the median scores between the two observations--the two distributions (e.g., pre-treatment and post-treatment) should display a similar shape when plotted as histograms or density plots.

If the shapes of the two distributions are markedly different (for instance, one is heavily skewed right and the other is symmetric), a significant WSRT result should be interpreted cautiously. In such cases, the significance indicates a difference between the distributions generally, but you cannot definitively attribute that difference solely to a shift in the central location (the median). Conversely, when the shapes are visually similar, a significant test allows the researcher to confidently state that the intervention or time difference resulted in a shift in the median value of the variable.

Guidelines for Utilizing the Wilcoxon Signed-Rank Test

Selecting the appropriate statistical test is critical for accurate analysis. The Wilcoxon Signed-Rank Test is the optimal choice when a specific set of research conditions and data characteristics are met. Researchers should proceed with the WSRT when their primary goal is to examine a difference between dependent measurements under non-parametric conditions, thereby ensuring the methodological choice aligns perfectly with the study design and data properties.

The decision matrix for employing the WSRT hinges on five key criteria. These criteria specify the nature of the research question, the type of data involved, and the sampling structure. If all the following conditions are satisfied, the Wilcoxon Signed-Rank Test is the most statistically sound method for the analysis:

You are seeking to establish a ****difference**** between the observations, not a relationship or

prediction.

Your measurements for the variable of interest must be **continuous** (or highly detailed ordinal).

The comparison involves **two and only two** related measurements or groups.

The samples are **paired** or dependent (repeated measures on the same subjects).

The distribution of the underlying variable is **skewed** or non-normal.

A clear understanding of these points helps to disambiguate the WSRT from other available statistical tests, allowing for correct application and interpretation.

Focus on Difference, Not Relationship

The Wilcoxon Signed-Rank Test is fundamentally a test of differences. It is employed when the research question asks whether a measurable change or disparity exists between the two paired observations. This contrasts sharply with other statistical objectives, such as determining the degree of association between two variables (which would require a correlation analysis) or building a model to forecast one variable based on others (which necessitates regression or prediction methods). When the goal is strictly to evaluate whether an intervention, treatment, or time lapse has produced a significant **shift** in scores, the WSRT is appropriate.

Requirement for Continuous Data

As previously established, the data must be continuous. This classification ensures that the mathematical operation of ranking the differences is valid. Continuous data permits fine granularity, covering measurements like reaction time, metabolic rate, precise financial income, or standardized psychological inventory scores. These measures allow for meaningful calculations of the difference score necessary for the test's ranking procedure.

Conversely, if the data is measured on a lower level scale, the WSRT should be avoided. Data types that are NOT continuous include ordered data (e.g., finishing rank in a marathon, severity ratings on a 1-5 scale), categorical data (e.g., political affiliation, geographic region, ethnicity), or binary data (e.g., success/failure, presence/absence of a disease). Using the WSRT on these variable types would violate the core assumption regarding the precision of the difference scores.

Limitation to Two Related Observations

A Wilcoxon Signed-Rank Test is structurally limited to comparing exactly two sets of measurements derived from a single group (Time 1 versus Time 2). The calculation of a single difference score for each participant mandates this binary comparison structure. This specificity is crucial in selecting the right non-parametric tool.

If you have three or more related observations (e.g., measuring performance at Baseline, Mid-

Treatment, and Post-Treatment), you must use a different analytical approach. If the variable is normally distributed across the time points, the One-Way Repeated Measures ANOVA is appropriate. If the variable remains skewed or non-normal across the multiple time points, the correct non-parametric extension is the Friedman Test.

The Requirement for Paired Samples (Dependence)

The WSRT is fundamentally a test for paired samples, meaning the two observations are linked, or dependent, often originating from the same subject. This dependency is typically established either through repeated measurements over time (longitudinal studies) or through matching pairs of participants based on external characteristics (matched-pairs designs). For example, if you randomly sample men at two points in time to get their IQ score, the two observations are paired because the pre-score is inherently linked to the post-score for that individual.

If you want to compare two independent groups, meaning your samples are completely separate and unrelated, then you probably want to use the Mann-Whitney U Test.

Confirming Non-Normality and Skewness

As discussed, the primary driver for choosing the WSRT over a paired t-test is the presence of non-normality or skewness in the data. Normality describes a symmetrical, bell-shaped distribution where the majority of data points cluster around the mean. When data deviates significantly from this shape, it is considered skewed, meaning it is leaning left or right with the majority of the data on one edge.

Researchers can visually inspect histograms or Q-Q plots to assess normality, but a more rigorous approach involves formal hypothesis testing. To objectively determine if your data violates the normality assumption, statistical tests such as the Kolmogorov-Smirnov test or the Shapiro-Wilk test should be employed. If these tests reject the null hypothesis of normality, the Wilcoxon Signed-Rank Test becomes the statistically safer choice for comparing the paired medians.

Practical Application Example: Evaluating an Intervention

To illustrate the practical utility of the Wilcoxon Signed-Rank Test, consider a health psychology study designed to evaluate the impact of a structured 12-week exercise program on physical fitness. The researchers measure the participants' maximal physical output at two distinct time points, creating a classic paired-samples design:

Observation 1: A group of people were evaluated at baseline.

Observation 2: This same group of people were evaluated after a 12-week exercise program.

Variable of interest: Number of pushups performed in 1 minute.

In this example, we have one group with two observations, meaning that the data are paired. Furthermore, physical performance scores are often highly skewed, meaning they are not normally distributed (skewed means leaning left or right with the majority of the data on the edge). This confirms the suitability of the WSRT for this analysis.

The null hypothesis (H_0), which is the statistical baseline scenario, posits that the exercise program has no effect, meaning there will be no difference in the median number of pushups performed before and after the program.

When we run the analysis, we obtain a test statistic (typically a Z or a T) and a p-value. The test statistic is a measure of how different the group's performance is on our pushups variable. The p-value is the probability of observing our results, or results more extreme, assuming the exercise program actually doesn't induce any change (i.e., assuming the null hypothesis is true). A p-value less than or equal to 0.05 means that our result is statistically significant and we can be confident that the difference observed is not due to chance alone.