

How to Perform a Mann-Whitney U Test in R: A Step-by-Step Guide

Authored by
stats writer

March 7, 2026

RECOMMENDED CITATION

stats writer (2026). *How to Perform a Mann-Whitney U Test in R: A Step-by-Step Guide*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=134470>

Introduction to Non-parametric Statistics and the Mann-Whitney U Test

In the vast landscape of quantitative analysis, the **Mann-Whitney U Test** stands as a cornerstone of **non-parametric statistical test** methodologies. Unlike parametric tests that rely on strict assumptions regarding the underlying distribution of data, such as normality, the Mann-Whitney U test offers a flexible alternative for researchers. It is specifically designed to compare the **median values** of two **independent** groups, making it an essential tool when the data is ordinal or when the requirements for a standard t-test cannot be met. By focusing on the ranks of the data points rather than their raw values, this test provides a resilient measure of central tendency differences that is less sensitive to outliers or extreme variance.

The utility of this test is particularly evident in fields such as clinical psychology, biology, and economics, where **distributions** are frequently skewed or sample sizes are inherently limited. When conducting research with small cohorts, the risk of violating the normality assumption is high, which can lead to inflated Type I errors or a loss of statistical power if parametric methods are misapplied. The Mann-Whitney U test, which is mathematically identical to the **Wilcoxon rank-sum test**, addresses these challenges by assessing whether one population tends to have larger values than the other. This makes it a robust choice for empirical studies where the goal is to determine if a specific intervention, such as a new medication or a policy change, has a **significant difference** on the outcome variable across two distinct populations.

In the **R** programming environment, performing this analysis is streamlined through built-in functions that handle the complex ranking and calculation of the **U statistic** and the associated **p-value**. The language provides a comprehensive suite of tools for data manipulation and statistical inference, allowing researchers to transition from data entry to result interpretation with high efficiency. Understanding how to correctly implement and interpret this test in **R** is a fundamental skill for any data scientist or statistician aiming to produce reliable, reproducible results. This guide will delve into the procedural nuances of the test, ensuring that users can confidently apply these methods to their own datasets while adhering to rigorous statistical standards.

Foundational Assumptions of the Mann-Whitney U Test

Before proceeding with any **statistical hypothesis testing**, it is imperative to verify that the data meets specific **assumptions** to ensure the validity of the results. For the **Mann-Whitney U Test**, the primary requirement is that the two samples being compared are **independent** of one another. This means that the observations in the first group must have no relationship or influence on the observations in the second group. For instance, if one group consists of patients receiving a **placebo** and another group consists of patients receiving a new drug, there should be no overlapping participants or matched pairs between these cohorts. If the data were paired, such as pre-test and post-test scores from the same individuals, a different test--the Wilcoxon signed-rank

test--would be required.

Another critical assumption involves the measurement scale of the data. The **Mann-Whitney U Test** requires that the dependent variable be at least ordinal, meaning that the data points can be ranked in a meaningful sequence. While the test is often used on continuous data that fails the normality check, it is equally applicable to discrete or ordinal data where the "distance" between values is not uniform. Furthermore, while the test does not require the data to follow a **normal distribution**, it does assume that the **distributions** of the two groups have a similar shape. If the shapes are similar, the test can be used to specifically compare **median values**; if the shapes differ significantly, the test instead evaluates whether one distribution stochastically dominates the other.

Finally, the test is most effective when the variance between the two groups is relatively consistent, although it is far more tolerant of heteroscedasticity than the **Student's t-test**. In practice, researchers often turn to this **non-parametric statistical test** when dealing with small **sample sizes** where the Central Limit Theorem cannot be relied upon to guarantee the normality of the sample mean. By acknowledging these foundational constraints, analysts can avoid common pitfalls and ensure that the **significant difference** they detect is truly reflective of the populations under study rather than an artifact of improper methodology.

The Statistical Mechanics Behind the Wilcoxon Rank-Sum Test

The internal logic of the **Mann-Whitney U Test** is based on a process of ranking all observations from both groups combined. Imagine pooling every data point from both samples and arranging them in ascending order. Each value is then assigned a rank: the smallest value receives a rank of 1, the second smallest a rank of 2, and so on. If there are tied values, they are typically assigned the average of the ranks they would have occupied. Once every observation has a rank, the ranks are summed for each individual group. The core idea is that if the two groups are drawn from the same population, the sum of their ranks should be relatively similar, adjusted for the respective size of each sample.

The calculation of the **U statistic** involves comparing these rank sums to the minimum possible sum that could occur by chance. Specifically, the test calculates how many times an observation from one group precedes an observation from the other group in the ordered list. This value, represented as "U," is then compared against a critical value from the Mann-Whitney distribution table or used to calculate a **p-value**. In **R**, this process is automated, but understanding the underlying ranking mechanism is vital for interpreting why certain data points--particularly those at the extremes--have such a profound impact on the final outcome of the test.

Because the test operates on ranks, it effectively "neutralizes" the impact of outliers that might otherwise skew the mean in a parametric test. For example, in a study measuring income, a single

billionaire would drastically increase the mean of a small sample, potentially leading to a false conclusion in a t-test. However, in a **non-parametric statistical test** like the Mann-Whitney, that billionaire is simply assigned the highest rank, the same as someone earning significantly less but still more than the rest of the group. This ranking approach ensures that the **median values** are the focal point of the comparison, providing a more accurate reflection of the "typical" experience within each group.

Practical Implementation Scenario: Panic Attack Intervention

To illustrate the application of the **Mann-Whitney U Test** in a real-world context, consider a clinical study involving a new pharmaceutical intervention for anxiety. Researchers are investigating whether a new drug can successfully reduce the frequency of panic attacks compared to a standard **placebo**. In this scenario, 12 patients are recruited and randomly assigned to one of two **independent** groups. Group A receives the new medication, while Group B receives the placebo. Over the course of a single month, each participant meticulously records the total number of panic attacks they experience, resulting in a dataset characterized by small **sample sizes** (n=6 per group) and likely non-normal counts.

The raw data for the month is as follows:

New Drug Group: 3, 5, 1, 4, 3, 5

Placebo Group: 4, 8, 6, 2, 1, 9

The primary objective is to determine if there is a statistically **significant difference** in the number of panic attacks between the two treatment arms. Given that the data consists of counts and the sample size is quite small, the **Mann-Whitney U Test** is the most appropriate choice for analysis. The researchers set a **level of significance** (alpha) of 0.05, meaning they are willing to accept a 5% risk of concluding a difference exists when it actually does not. By utilizing **R**, the researchers can quickly process these numbers to find the W statistic (which **R** uses to represent the rank-sum) and the corresponding **p-value**.

The following table summarizes the experimental results for better visualization before we move into the **R** code implementation:

NEW DRUG	PLACEBO
3	4
5	8
1	6
4	2

NEW DRUG	PLACEBO
3	1
5	9

Methodology One: Executing the Test Using Independent Vectors

The first and perhaps most straightforward way to perform the **Mann-Whitney U Test** in **R** is by organizing the data into two separate numeric vectors. This method is highly efficient for quick comparisons or when data is provided in a simple list format. To begin, we define the "new" and "placebo" groups as vectors using the `c()` function. Once the data is stored in these variables, we invoke the `wilcox.test` function, passing the two vectors as arguments. This function is the standard tool in **R** for performing the Wilcoxon rank-sum analysis, which serves as the implementation of the Mann-Whitney U logic.

By default, `wilcox.test` performs a two-sided test, which evaluates the **null hypothesis** that there is no difference in the distribution of panic attacks between the two groups. The function also applies a **continuity correction** by default, which is a technical adjustment used to improve the accuracy of the **p-value** when approximating a discrete distribution with a continuous one. This is particularly useful for small datasets like the one in our example. The output provides the *W* statistic and the p-value, which are the essential components for our conclusion.

Below is the **R** code and the resulting output for this vector-based approach:

```
#create a vector for each group
new <- c(3, 5, 1, 4, 3, 5)
placebo <- c(4, 8, 6, 2, 1, 9)

#perform the Mann Whitney U test
wilcox.test(new, placebo)

#output
Wilcoxon rank sum test with continuity correction

data: new and placebo
W = 13, p-value = 0.468
alternative hypothesis: true location shift is not equal to 0
```

In this output, the value of $W = 13$ represents the rank sum calculated by the algorithm. The **p-value** of 0.468 is the critical value we must analyze against our alpha of 0.05. Because 0.468 is significantly higher than 0.05, the data suggests that any observed difference between the drug

group and the **placebo** group is likely due to chance rather than the drug's effectiveness. This initial step provides a clear mathematical basis for accepting the **null hypothesis** and suggests that more research or a larger sample might be necessary.

Methodology Two: Utilizing Data Frames for Formula-Based Testing

In more complex data analysis workflows, it is common to store information within a **data frame** rather than separate vectors. This structure is often preferred because it allows for better data management, especially when dealing with multiple variables or larger datasets imported from CSV files or databases. In this approach, one column represents the dependent variable (the number of panic attacks), while a second column serves as a grouping factor (identifying whether the participant received the drug or the **placebo**). This "tidy" data format is standard in modern data science and works seamlessly with **R**'s formula syntax.

To implement the **Mann-Whitney U Test** using this method, we utilize the formula notation ``response ~ group``. This tells **R** to analyze the "attacks" variable based on the categories found in the "drug_group" column. The **wilcox.test** function is versatile enough to recognize this syntax, provided we specify the data frame in the ``data`` argument. This methodology is particularly useful when you are working with the "tidyverse" suite of packages or when you need to perform the test across multiple subsets of data using functions like ``lapply`` or ``group_by``.

The following code block demonstrates how to construct the **data frame** and execute the test using the formula method:

```
#create a data frame with two columns, one for each group  
drug_data <- data.frame(attacks = c(3, 5, 1, 4, 3, 5, 4, 8, 6, 2, 1, 9),  
drug_group = c(rep("old", 6), rep("placebo", 6)))
```

```
#perform the Mann Whitney U test  
wilcox.test(attacks~drug_group, data = drug_data)
```

```
#output  
data: attacks by drug_group  
W = 13, p-value = 0.468  
alternative hypothesis: true location shift is not equal to 0
```

As expected, the results are identical to the vector-based method, yielding a *W* statistic of 13 and a **p-value** of 0.468. This consistency reinforces the fact that the choice between vectors and data frames is purely one of organizational preference and workflow compatibility. Whether you are performing a quick ad-hoc check or building a robust analysis pipeline, **R** ensures that the underlying **non-parametric statistical test** remains accurate and reliable.

Interpreting the Output: Analyzing P-Values and Test Statistics

Interpreting the results of a **Mann-Whitney U Test** requires a clear understanding of the **null hypothesis** (H0) and the **alternative hypothesis** (H1). In our panic attack study, the null hypothesis posits that there is no difference in the distribution of panic attack counts between the drug group and the **placebo** group. Conversely, the alternative hypothesis suggests that a difference does exist--that one group tends to have higher or lower values than the other. The **p-value** is the probability of observing a result at least as extreme as the one obtained, assuming the null hypothesis is true.

With a **p-value** of 0.468, we find that our result is well above the traditional **level of significance** threshold of 0.05. Consequently, we fail to reject the null hypothesis. In practical terms, this means that the study does not provide sufficient evidence to claim that the new drug is effective at reducing panic attacks compared to the placebo. It is important to note that "failing to reject" the null hypothesis is not the same as "proving" the null hypothesis is true; it simply means the current data is not strong enough to conclude otherwise. This distinction is vital in **statistical hypothesis testing**.

The W statistic itself, which is 13 in our case, represents the sum of the ranks for the first group minus the minimum possible rank sum. While the p-value is usually the primary focus for decision-making, the W value can be useful for calculating effect sizes, such as the Rank-Biserial Correlation. An effect size provides a measure of the magnitude of the difference, which can be helpful even when the **significant difference** is not reached. For instance, a small p-value in a very large sample might indicate a statistically significant but practically meaningless difference, whereas a large p-value in a small sample might mask a potentially important trend that warrants further investigation with more participants.

Alternative Hypotheses and One-Tailed Testing Procedures

In many research scenarios, the investigator has a specific expectation about the direction of the effect. For example, if we are testing a new drug, we might not just be looking for "any" difference; we specifically want to know if the drug results in *fewer* panic attacks than the **placebo**. This is known as a one-tailed test or a directional **alternative hypothesis**. In **R**, this is easily accommodated by adding the ``alternative`` parameter to the **wilcox.test** function. By specifying ``alternative = "less"``, we are testing if the values in the first group are stochastically smaller than the values in the second group.

When we switch to a one-tailed test, the **p-value** is essentially halved if the observed difference is in the predicted direction. In our example, the two-sided p-value was 0.468. When we re-run the test with the "less" alternative, the p-value becomes 0.234. While this is a more "favorable" number

for the researcher, it still remains well above the 0.05 **level of significance**. This means that even with a more specific hypothesis, the evidence remains insufficient to support the drug's efficacy. Using directional tests should always be justified by prior theory or logical reasoning, as it effectively makes it "easier" to find significance, which can increase the risk of false positives if used inappropriately.

The following code demonstrates the implementation of a one-tailed **Mann-Whitney U Test** in **R**:

```
#create a vector for each group
```

```
new <- c(3, 5, 1, 4, 3, 5)
```

```
placebo <- c(4, 8, 6, 2, 1, 9)
```

```
#perform the Mann Whitney U test, specify alternative="less"
```

```
wilcox.test(new, placebo, alternative="less")
```

```
#output
```

```
Wilcoxon rank sum test with continuity correction
```

```
data: new and placebo
```

```
W = 13, p-value = 0.234
```

```
alternative hypothesis: true location shift is less than 0
```

The **U statistic** (W) remains 13, but the **p-value** adjustment reflects the directional nature of the inquiry. This flexibility highlights why **R** is such a powerful environment for **non-parametric statistical test** procedures. It allows the user to tailor the analysis precisely to their research questions while maintaining the integrity of the underlying mathematical model. Analysts should always report whether they used a one-tailed or two-tailed test in their final documentation to ensure transparency and reproducibility in their findings.

Conclusion: Comparative Analysis and Best Practices

The **Mann-Whitney U Test** serves as a vital tool in the statistician's arsenal, particularly when the rigid **assumptions** of parametric methods cannot be satisfied. By focusing on ranks rather than raw values, it provides a reliable way to compare **median values** between **independent** groups without the need for a **normal distribution**. While the **Student's t-test** is more powerful when data is truly normal, the Mann-Whitney U test is often nearly as powerful and far more robust when data is skewed or contains outliers. Choosing the right test is a matter of understanding your data's structure and the specific questions you aim to answer.

When implementing these tests in **R**, researchers benefit from the high-level **wilcox.test** function, which handles both vector-based and formula-based inputs. This versatility, combined with options

for **continuity correction** and directional **alternative hypothesis** specifications, makes **R** an industry-standard choice for non-parametric analysis. Best practices suggest always performing a visual inspection of your data--using boxplots or histograms--before running the test to check for similar distribution shapes and to identify any potential outliers that might influence the results.

In summary, while our specific example did not yield a **significant difference**, the process demonstrated the efficiency and clarity of the **Mann-Whitney U Test**. By adhering to the proper **assumptions** and utilizing the robust features of **R**, analysts can draw meaningful conclusions about their data with a high degree of confidence. Whether you are comparing the effectiveness of medical treatments, the results of different marketing strategies, or environmental changes across locations, the Mann-Whitney U test remains a reliable, non-parametric foundation for **statistical hypothesis testing**.

ARABPSYCHOLOGY.COM