

How to Perform a Mann-Whitney U Test to Compare Two Independent Samples

Authored by
stats writer

February 28, 2026

RECOMMENDED CITATION

stats writer (2026). *How to Perform a Mann-Whitney U Test to Compare Two Independent Samples*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=133251>

Introduction to Non-Parametric Comparison

In the field of **inferential statistics**, researchers frequently encounter scenarios where they must determine if two separate groups differ significantly from one another. While the **independent samples t-test** is often the first choice for such comparisons, it relies heavily on the assumption that the data follows a **normal distribution**. When data is skewed, contains outliers, or is measured on an **ordinal scale**, these parametric assumptions are violated, rendering the t-test unreliable. In such instances, the **Mann-Whitney U Test** serves as a powerful and robust alternative, allowing for a rigorous comparison of two **independent samples** without requiring the stringent parameters of normality.

The **Mann-Whitney U Test**, which is also widely recognized as the **Wilcoxon rank-sum test**, is a **nonparametric** procedure designed to assess whether two populations have the same distribution. Unlike parametric tests that focus on the difference between means, this test evaluates the **ranks** of the data points. By converting raw scores into ranks, the test effectively minimizes the influence of extreme **outliers** and focuses instead on the relative positioning of values within the combined dataset. This makes it an essential tool for researchers working with small **sample sizes** or non-interval data where the mean might not be a representative measure of central tendency.

Fundamentally, the test operates on the principle of **stochastic dominance**. It seeks to determine the probability that a randomly selected observation from one population will be greater than a randomly selected observation from the second population. If the two samples come from the same distribution, the **U statistic** will reflect a balanced distribution of ranks across both groups. Conversely, if one group consistently ranks higher than the other, the resulting **U statistic** will deviate significantly from the expected value under the **null hypothesis**, suggesting that the distributions are indeed distinct. This conceptual framework allows the **Mann-Whitney U Test** to provide high **statistical power** even when the underlying population parameters are unknown or complex.

Identifying the Mann-Whitney U Test and Its Applications

The primary utility of the **Mann-Whitney U Test** lies in its flexibility across various scientific disciplines, particularly when the criteria for a **t-test** cannot be met. It is specifically recommended when **sample sizes** are small--typically defined as having fewer than 30 observations per group--and when the **sampling distribution** of the means cannot be assumed to be normal. Because it does not assume a specific bell-shaped curve for the population, it is highly effective for analyzing **real-world data** that often presents as asymmetrical or heavily tailed. By utilizing the **ranks** of the data, the test remains valid even when the data is strictly **ordinal**, such as survey responses or performance ratings.

Consider a scenario in **labor economics** where a researcher wishes to compare the starting salaries of five graduates from University A against five graduates from University B. In such a small dataset, a single high-earning individual could significantly inflate the mean, leading to a type I error in a parametric test. However, by applying the **Mann-Whitney U Test**, the researcher focuses on the **rank order** of the salaries. If University A's graduates consistently occupy the highest rank positions, the test will detect a statistically significant difference regardless of the exact numerical gap between the salaries, providing a more reliable conclusion about the economic outcomes of the two institutions.

Similarly, the test is indispensable in **health sciences** and **clinical research**. For instance, if a nutritionist is measuring weight loss across two small groups of participants using different diets, the **weight loss** measurements might not follow a **normal distribution**. Using the **Wilcoxon rank-sum test** allows the researcher to compare the efficacy of Diet A versus Diet B without being misled by individual variations in metabolism that might skew the average. Whether the research involves comparing **exam scores** between two small classrooms or assessing the recovery times of patients under different therapeutic protocols, this test ensures that the lack of normality does not hinder the ability to draw meaningful statistical inferences.

Essential Assumptions for Statistical Validity

While the **nonparametric** nature of the **Mann-Whitney U Test** provides significant freedom from the assumption of normality, it is not entirely free of requirements. To ensure that the results are valid and the **p-value** is accurate, researchers must verify that four primary **assumptions** are satisfied. The first assumption concerns the **measurement scale** of the **dependent variable**. The data must be at least **ordinal** or **continuous**. **Ordinal** variables, such as **Likert items** ranging from "Strongly Disagree" to "Strongly Agree," allow for ranking, while **continuous** variables like height or weight provide a precise numerical sequence that can be converted into ranks.

The second critical assumption is the **independence of observations**. This means that there should be no relationship between the observations in each group or between the groups themselves. For the **Mann-Whitney U Test** to be appropriate, the participants in Group A must be different from those in Group B, and the data point of one participant should not influence the data point of another. If the samples were related--such as measuring the same individual before and after a treatment--the **Wilcoxon Signed-Rank Test** would be required instead. **Independence** ensures that the calculated **U statistic** accurately reflects the differences between two distinct populations rather than internal correlations within the data.

The third assumption involves the **shape of the distribution**. While the test does not require the distributions to be normal, it does assume that the **shapes** of the distributions for both groups are relatively similar. If the two distributions have the same shape, the **Mann-Whitney U Test** can be

used to specifically determine if there is a difference in the **medians** of the two groups. However, if the shapes differ significantly, the test is instead interpreted as assessing whether there is a general difference in the distributions (i.e., whether one group tends to have higher values than the other). Finally, **random sampling** is assumed to ensure that the findings can be generalized to the broader populations from which the samples were drawn.

Deciphering the Mathematical Logic of Ranking

The core of the **Mann-Whitney U Test** is the process of **ranking** the combined data from both groups. To calculate the **test statistic**, all observations from Sample 1 and Sample 2 are pooled together and arranged in ascending order. Each value is then assigned a rank, starting with 1 for the smallest value and continuing through N (the total number of observations). In cases where **ties** occur--meaning two or more observations have the exact same value--the average of the ranks they would have otherwise occupied is assigned to each. For example, if the 4th and 5th values are identical, both are assigned a rank of 4.5. This conversion from raw data to **ranks** is what allows the test to be **nonparametric** and resistant to the influence of extreme **outliers**.

Once the ranks are assigned, they are summed for each group separately, resulting in the values **R1** and **R2**. These sums are then used to calculate two separate U values, **U1** and **U2**, using the following formulas:

$$U1 = n1 * n2 + - R1$$

$$U2 = n1 * n2 + - R2$$

In these equations, **n1** and **n2** represent the **sample sizes** for each respective group. The term represents the sum of ranks for a group if all its observations were the smallest in the combined set. By subtracting the actual **sum of ranks** (R), the formula determines how many times an observation from one group "wins" or precedes an observation from the other group. The final **U statistic** used for **hypothesis testing** is simply the smaller of the two values (U1 or U2). A smaller U value indicates a more significant separation between the two groups, as it suggests that one group consistently outranked the other.

It is important to note that the maximum possible value for U is the product of the two sample sizes (**n1 * n2**). If the two samples are perfectly separated--meaning all values in one group are smaller than every value in the other--one of the U values will be zero. Conversely, if the ranks are perfectly interspersed, the U value will be approximately half of **n1 * n2**. This mathematical balance is what allows the test to effectively measure the overlap between two **independent samples**. By comparing this calculated value to a **critical value** from a **Mann-Whitney U distribution table**, researchers can determine the probability that the observed difference occurred by **random chance**.

Comprehensive Methodology for Hypothesis Testing

Executing a **Mann-Whitney U Test** follows a rigorous five-step **hypothesis testing** procedure common in **quantitative analysis**. The first step is to formally state the **null hypothesis** (H_0) and the **alternative hypothesis** (H_a). Typically, the **null hypothesis** posits that there is no difference between the two populations, implying that a randomly selected value from the first group is equally likely to be larger or smaller than a randomly selected value from the second. The **alternative hypothesis**, usually two-sided, suggests that the two populations are not equal and that one distribution is stochastically dominant over the other.

The second step involves selecting a **significance level** (α), which is the threshold for rejecting the **null hypothesis**. In most scientific research, an α of 0.05 is used, though more stringent levels like 0.01 may be applied in clinical or high-stakes environments. The third step is the calculation of the **test statistic** (U). As previously detailed, this involves ranking the combined data, summing the ranks for each sample, and applying the **U-test formulas**. This phase requires meticulous attention to detail, especially when dealing with **ties** in the data, as inaccurate ranking will lead to an incorrect **U statistic** and potentially a faulty conclusion.

In the fourth step, the researcher must decide whether to **reject or fail to reject the null hypothesis**. This is done by comparing the calculated U value to a **critical value** found in a standardized reference table. Unlike the **t-test** or **F-test**, where a larger statistic indicates significance, the **Mann-Whitney U Test** is significant if the calculated U is **less than or equal to the critical value**. If the U value is small enough, it provides sufficient evidence to conclude that the difference between the groups is unlikely to have arisen from **sampling error** alone. Finally, the fifth step is the **interpretation of results**, where the mathematical finding is translated back into the context of the original research question, providing a clear answer to the scientific inquiry.

Practical Application: Evaluating Clinical Drug Efficacy

To illustrate the application of the **Mann-Whitney U Test** in a clinical setting, consider a study designed to evaluate the effectiveness of a new pharmacological treatment for **panic attacks**. In this experiment, 12 patients are divided into two equal groups: one receiving a new drug and the other receiving a **placebo**. Over the course of one month, the number of panic attacks experienced by each patient is recorded. Because the sample size is very small ($n=6$ per group) and the number of attacks is unlikely to be **normally distributed**, the **nonparametric U-test** is the most appropriate analytical tool.

The raw data for the two groups is presented in the table below:

NEW DRUG	PLACEBO
3	4
5	8
1	6
4	2
3	1
5	9

In accordance with the **hypothesis testing** framework, we set our **null hypothesis** as the equality of both groups, with an alpha level of 0.05. To find the **U statistic**, we combine and rank the 12 observations: 1, 1, 2, 3, 3, 4, 4, 5, 5, 6, 8, 9. The resulting **ranks** are 1.5, 1.5, 3, 4.5, 4.5, 6.5, 6.5, 8.5, 8.5, 10, 11, 12. Summing these ranks for each group gives us **R1 = 34** for the drug group and **R2 = 44** for the placebo group. Applying the formulas: **U1 = (6*6) + - 34 = 23**, and **U2 = (6*6) + - 44 = 13**. Our final test statistic is **U = 13**.

To determine significance, we consult the **critical value** table for n1=6 and n2=6 at a 0.05 significance level. The table indicates a **critical value** of 5.

n1 \ n2	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
2							0	0	0	0	1	1	1	1	1	2	2	2	2
3				0	1	1	2	2	3	3	4	4	5	5	6	6	7	7	8
4			0	1	2	3	4	4	5	6	7	8	9	10	11	11	12	13	14
5		0	1	2	3	5	6	7	8	9	11	12	13	14	15	17	18	19	20
6		1	2	3	5	6	7	10	11	13	14	16	17	19	21	22	24	25	27
7		1	3	5	6	8	10	12	14	16	18	20	22	24	26	28	30	32	34
8	0	2	4	6	7	10	13	15	17	19	22	24	26	29	31	34	36	38	41
9	0	2	4	7	10	12	15	17	20	23	26	28	31	34	37	39	42	45	48
10	0	3	5	8	11	14	17	20	23	26	29	33	36	39	42	45	48	52	55
11	0	3	6	9	13	16	19	23	26	30	33	37	40	44	47	51	55	58	62
12	1	4	7	11	14	18	22	26	29	33	37	41	45	49	53	57	61	65	69
13	1	4	8	12	16	20	24	28	33	37	41	45	50	54	59	63	67	72	76
14	1	5	9	13	17	22	26	31	36	40	45	50	55	59	64	67	74	78	83
15	1	5	10	14	19	24	29	34	39	44	49	54	59	64	70	75	80	85	90
16	1	6	11	15	21	26	31	37	42	47	53	59	64	70	75	81	86	92	98
17	2	6	11	17	22	28	34	39	45	51	57	63	67	75	81	87	93	99	105
18	2	7	12	18	24	30	36	42	48	55	61	67	74	80	86	93	99	106	112
19	2	7	13	19	25	32	38	45	52	58	65	72	78	85	92	99	106	113	119
20	2	8	13	20	27	34	41	48	55	62	69	76	83	90	98	105	112	119	127

Because our calculated **U statistic** of 13 is greater than the **critical value** of 5, we **fail to reject the null hypothesis**. Statistically, this means that the evidence is insufficient to claim a significant difference in the frequency of panic attacks between those taking the new drug and those taking the **placebo**. In a clinical context, this might suggest that while the drug could have some effect, the current study lacks the **statistical power** or the treatment lacks the potency to demonstrate a clear advantage over the placebo within this specific sample.

Empirical Analysis: Assessing Educational Study Habits

Another compelling example of the **Mann-Whitney U Test** involves educational research. Suppose a researcher wants to investigate whether a specific study regimen--30 minutes of daily review--improves **test scores** compared to a control group with no specific study requirements. A total of 15 students are assigned to either the "Study" group (n=8) or the "No-Study" group (n=7). Because **academic performance** data in small groups is often skewed or contains outliers, a **nonparametric** approach is chosen to ensure the integrity of the findings.

The test scores for the two groups are listed in the following table:

STUDY	NO-STUDY
89	88
92	93
94	95
96	75
91	72
99	80
84	81
90	

Following the same **ranking procedure**, all 15 scores are combined and ranked from 1 to 15. The sum of ranks for the Study group is calculated as **R1 = 80**, while the sum of ranks for the No-Study group is **R2 = 40**. Using the **U-test formulas** with n1=8 and n2=7: **U1 = (8*7) + - 80 = 12**, and **U2 = (8*7) + - 40 = 44**. The smaller value, **U = 12**, is our **test statistic**. For this analysis, a more conservative **significance level** of 0.01 is selected to minimize the risk of a false positive.

Consulting the **critical value** table for n1=8 and n2=7 at the 0.01 level, we find a **critical value** of 6.

n1 \ n2	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
2																			
3								0	0	0	1	1	1	2	2	2	2	3	3
4					0	0	1	1	2	2	3	3	4	5	5	6	6	7	8
5				0	1	1	2	3	4	5	6	7	7	8	9	10	11	12	13
6			0	1	2	3	4	5	6	7	9	10	11	12	13	15	16	17	18
7			0	1	3	4	6	7	9	10	12	13	15	16	18	19	21	22	24
8			1	2	4	6	7	9	11	13	15	17	18	20	22	24	26	28	30
9		0	1	3	5	7	9	11	13	16	18	20	22	24	27	29	31	33	36
10		0	2	4	6	9	11	13	16	18	21	24	26	29	31	34	37	39	42
11		0	2	5	7	10	13	16	18	21	24	27	30	33	36	39	42	45	46
12		1	3	6	9	12	15	18	21	24	27	31	34	37	41	44	47	51	54
13		1	3	7	10	13	17	20	24	27	31	34	38	42	45	49	53	56	60
14		1	4	7	11	15	18	22	26	30	34	38	42	46	50	54	58	63	67
15		2	5	8	12	16	20	24	29	33	37	42	46	51	55	60	64	69	73
16		2	5	9	13	18	22	27	31	36	41	45	50	55	60	65	70	74	79
17		2	6	10	15	19	24	29	34	39	44	49	54	60	65	70	75	81	86
18		2	6	11	16	21	26	31	37	42	47	53	58	64	70	75	81	87	92
19	0	3	7	12	17	22	28	33	39	45	51	56	63	69	74	81	87	93	99
20	0	3	8	13	18	24	30	36	42	46	54	60	67	73	79	86	92	99	105

Since our calculated **U statistic** (12) remains higher than the **critical value** (6), we again **fail to reject the null hypothesis**. Despite the Study group having a higher rank sum, the difference is not pronounced enough to be statistically significant at the 0.01 level. This result indicates that, within the constraints of this small sample, there is no definitive evidence that the 30-minute study intervention lead to superior **test scores**. This highlights the importance of **sample size** in detecting subtle effects and the rigor provided by the **Mann-Whitney U Test** in preventing over-interpretation of raw data differences.

The Role of Critical Values and Statistical Significance

The interpretation of the **Mann-Whitney U Test** hinges on the relationship between the **U statistic** and the **critical value**. In many other statistical tests, such as the **Z-test** or **Chi-Square test**, a larger **test statistic** corresponds to a smaller **p-value** and a higher likelihood of significance. However, the **U-test** operates inversely. A smaller U value represents less overlap between the ranks of the two groups, which in turn indicates a more significant difference. When the U value is zero, it means the two samples are completely distinct with no overlapping values, providing the strongest possible evidence against the **null hypothesis**.

Finding the correct **critical value** requires a specialized table that accounts for the sizes of both samples (n1 and n2) and the chosen alpha level. For **sample sizes** larger than 20, the distribution

of the **U statistic** begins to approximate a **normal distribution**. In these cases, researchers often use a **Z-score** transformation to calculate the **p-value** directly. This **normal approximation** allows the **Mann-Whitney U Test** to be scaled for larger datasets while maintaining its **nonparametric** advantages. However, for the small samples typically associated with this test, the exact **critical value** approach remains the standard for accuracy.

It is also vital to distinguish between **one-tailed** and **two-tailed** tests when looking up **critical values**. A **two-tailed test**, which is the most common, assesses whether there is any difference between the groups regardless of direction. A **one-tailed test** is used only when a researcher has a specific, pre-defined hypothesis that one group will perform better than the other. Choosing the wrong tail can lead to incorrect conclusions about **statistical significance**. By strictly adhering to these protocols, the **Mann-Whitney U Test** provides a reliable mechanism for researchers to validate their findings and ensure that their conclusions are supported by the mathematical reality of the data ranks.

Distinguishing the Mann-Whitney U Test from Parametric Alternatives

The decision to use the **Mann-Whitney U Test** over the **independent samples t-test** is a fundamental choice in **experimental design**. The **t-test** is generally more powerful when the data perfectly meets the assumption of **normality** and **homoscedasticity** (equal variances). However, the **t-test's** reliance on the mean makes it highly sensitive to **outliers**, which can lead to a **Type II error** if the outlier masks a true difference, or a **Type I error** if the outlier creates a false one. The **Mann-Whitney U Test**, by using **median-based ranking**, provides a safeguard against such anomalies, ensuring that the results reflect the overall trend of the data rather than the influence of a few extreme points.

Furthermore, the **Mann-Whitney U Test** is often preferred in the **social sciences** where data is frequently collected via subjective scales. Since these scales do not have a consistent "distance" between points (e.g., the difference between "Happy" and "Very Happy" may not be the same as between "Neutral" and "Happy"), they are **ordinal** rather than **interval**. Parametric tests like the **t-test** are technically inappropriate for **ordinal data**, whereas the **Wilcoxon rank-sum test** handles it perfectly. This makes the **U-test** a more versatile tool for researchers who must navigate the complexities of human behavior and opinion.

In summary, the **Mann-Whitney U Test** is an essential component of the statistical toolkit, offering a robust method for comparing **independent samples** when parametric conditions fail. Its focus on **ranks** rather than means provides a unique perspective on data distribution and stochastic dominance. Whether used in clinical trials, educational assessments, or economic comparisons, it ensures that researchers can make confident, data-driven decisions. By understanding the **assumptions**, mathematical formulas, and **hypothesis testing** steps involved, one can effectively

leverage this test to uncover significant insights within even the most challenging datasets.

Additional Resources and Practical Implementation

For those interested in the practical application of these statistical methods in modern programming environments, there are numerous resources available. Mastering the **Mann-Whitney U Test** is not just about manual calculation; it is also about leveraging technology to perform these tests on larger datasets efficiently. Many researchers now utilize **Python** or **R** to automate the ranking and U-calculation processes, which significantly reduces the risk of human error and allows for more complex **data visualization**.

How to Perform a Mann-Whitney U Test in Python

Understanding the nuances of **nonparametric statistics** opens the door to a more sophisticated level of **data analysis**. As data becomes increasingly complex and non-traditional, the reliance on robust tests like the **Mann-Whitney U Test** will only continue to grow. By prioritizing high-quality, **peer-reviewed sources** and official documentation, students and professionals alike can ensure they are applying these methods correctly and contributing valid, reliable findings to their respective fields of study.