

Parametric and Non-Parametric Tests

Every hypothesis test makes assumptions. A z-test for a mean assumes the population is normal, or that n is large enough for the Central Limit Theorem to rescue us. A test on a correlation coefficient ρ assumes the underlying data is bivariate normal. Tests like these, which assume the population has a particular distributional form and then test a *parameter* of it, are called **parametric tests**.

What if the assumptions fail? With a sample of size 8 from a skewed, unknown distribution, we cannot assume normality and the CLT is no help. We need tests that work regardless.

Definition. A **non-parametric test** is a hypothesis test that makes few or no assumptions about the distribution of the underlying population.

Non-parametric tests are useful when:

- the sample is small and the population cannot be assumed normal;
- the data is clearly skewed or contains outliers (medians behave better than means);
- the data is only *ordinal* — ranks or preferences rather than true measurements.

The price: when parametric assumptions *do* hold, non-parametric tests are less powerful (they throw away information, as we shall see).

Remark. You have already met one non-parametric test. Testing $H_0: \rho = 0$ with the product-moment correlation coefficient requires bivariate normality; testing with Spearman's rank correlation coefficient r_s does not, because ranking the data discards its distribution. See the Correlation notes.

Because we drop all distributional assumptions, our hypotheses change character: non-parametric tests are about the **median** m rather than the mean (the median always exists and is meaningful for any shape of distribution).

Textbook Exercises: [CUPS] Ch 4 §1; [S3&4] S4 Ch 2

The Single-Sample Sign Test

The simplest possible test. To test $H_0: m = m_0$: if the population median really is m_0 , then each observation is equally likely to fall above or below m_0 . So the number of observations above m_0 is binomial with $p = \frac{1}{2}$.

- Tip (The procedure)** 1. $H_0: m = m_0$; $H_1: m \neq m_0$ (or $>$, $<$), where m is the population median of [context].
2. **Discard** any observations exactly equal to m_0 . Let n be the number remaining.
 3. Record the sign (+ or -) of each remaining observation minus m_0 , and count the positive signs: $S \sim B(n, \frac{1}{2})$ under H_0 .
 4. Compute the tail probability of a result at least as extreme as the one observed; for a two-tailed test, double it. Compare with the significance level.

Example (Sign test)

A manufacturer claims the median lifetime of its batteries is more than 30 hours. Twelve batteries are tested, with lifetimes (hours):

31, 34, 28, 36, 30, 32, 41, 29, 33, 38, 35, 31.

Test the claim at the 5% significance level using a sign test.

$H_0: m = 30$; $H_1: m > 30$, where m is the population median battery lifetime in hours.
 One observation equals 30: discard it, leaving $n = 11$. Signs of (lifetime - 30):

+, +, -, +, +, +, -, +, +, +, +

so 9 positive signs out of 11. Under H_0 , the number of positive signs $S \sim B(11, \frac{1}{2})$.

$$P(S \geq 9) = \frac{\binom{11}{9} + \binom{11}{10} + \binom{11}{11}}{2^{11}} = \frac{55 + 11 + 1}{2048} = \frac{67}{2048} = 0.0327$$

Since $0.0327 < 0.05$, we reject H_0 . There is evidence at the 5% level to suggest that the median battery lifetime exceeds 30 hours.

Example (OCR S4, June 2012)

A one-tail sign test of a population median is to be carried out at the 5% significance level using a sample of size n .

- (i) Show by calculation that the test can never result in rejection of the null hypothesis when $n = 4$.
- (ii) The coach of a college swimming team expects Elena, the best 50m freestyle swimmer, to have a median time less than 30 seconds. Elena found from records of her previous 72 swims that 44 were less than 30 seconds and 28 were greater than 30 seconds. Stating a necessary assumption, test at the 5% significance level whether Elena's median time for the 50m freestyle is less than 30 seconds.

(i) The most extreme one-tailed result possible is all four signs the same, with probability

$$\left(\frac{1}{2}\right)^4 = \frac{1}{16} = 0.0625 > 0.05,$$

so even the most extreme possible sample is not significant at the 5% level: H_0 can never be rejected when $n = 4$.

(ii) Assumption: the 72 recorded swims are a random sample of Elena's times. $H_0: m = 30$; $H_1: m < 30$, where m is the population median of Elena's 50 m times.

Let S be the number of swims under 30 seconds; under H_0 , $S \sim B(72, \frac{1}{2})$, and we observe $S = 44$. Since $n = 72$ is large, $S \approx N(36, 18)$, and with a continuity correction

$$\mathbb{P}(S \geq 44) \approx \mathbb{P}\left(Z > \frac{43.5 - 36}{\sqrt{18}}\right) = \mathbb{P}(Z > 1.768) = 0.0385.$$

Since $0.0385 < 0.05$, we reject H_0 . There is evidence at the 5% level to suggest that Elena's median 50 m freestyle time is less than 30 seconds.

Remark (The weakness of the sign test). The sign test uses almost no assumptions — but also almost no information. An observation 0.1 above the median and one 40 above the median count identically. Data like $-1, -2, -1, +30, +45, +38, +29$ would be treated as “3 minus, 4 plus”, ignoring the obvious message in the magnitudes. The next test repairs this.

Textbook Exercises: [CUP.S] Ch 4 §1; [S3&4] S4 Ch 2

The Wilcoxon Signed-Rank Test

The Wilcoxon signed-rank test uses the *sizes* of the deviations from m_0 as well as their signs — at the cost of one extra assumption.

Fact (Assumption) — The Wilcoxon signed-rank test assumes the population distribution is **symmetric**. Under H_0 , deviations of any given size are then equally likely to be positive or negative.

- Tip (The procedure)**
1. $H_0: m = m_0$; $H_1: m \neq m_0$ (or one-tailed).
 2. Compute the differences $d_i = x_i - m_0$, discarding any zeros.
 3. Rank the $|d_i|$ from 1 (smallest) to n (largest). (We do not consider tied ranks.)
 4. Let T^+ = sum of the ranks of the positive differences, T^- = sum of ranks of the negative differences. (Also written W^+ and W^- .) Check: $T^+ + T^- = \frac{n(n+1)}{2}$.
 5. Test statistic $T = \min(T^+, T^-)$. Reject H_0 if $T \leq$ the critical value from the formula-booklet table (small T means one side's ranks are suspiciously light).

Example (Single-sample Wilcoxon signed-rank test)

Bags of flour are sold as having median weight 50g. A sample of 10 bags weighs (g):

52, 43, 38, 45, 41, 35, 51, 40, 46, 42.

Assuming the distribution of weights is symmetric, test at the 5% significance level whether the population median weight differs from 50g.

$H_0: m = 50$; $H_1: m \neq 50$, where m is the population median weight.

x_i	52	43	38	45	41	35	51	40	46	42
$d_i = x_i - 50$	+2	-7	-12	-5	-9	-15	+1	-10	-4	-8
rank of $ d_i $	2	5	9	4	7	10	1	8	3	6

$$T^+ = 2 + 1 = 3, \quad T^- = 5 + 9 + 4 + 7 + 10 + 8 + 3 + 6 = 52$$

(Check: $3 + 52 = 55 = \frac{10 \times 11}{2}$.) So $T = \min(T^+, T^-) = 3$.

From the booklet, the critical value for $n = 10$ at the 5% two-tailed level is 8.

Since $T = 3 \leq 8$, we reject H_0 . There is evidence at the 5% level to suggest that the population median weight of the bags is not 50g (the sample suggests it is lower).

Fact — Under H_0 ,

$$\mathbb{E}[T^+] = \mathbb{E}[T^-] = \frac{n(n+1)}{4}.$$

Values of T far below $\frac{n(n+1)}{4}$ are evidence against H_0 .

The proof is one line.

Under H_0 , each rank i joins T^+ independently with probability $\frac{1}{2}$ (by symmetry of the distribution about m_0), so

$$\mathbb{E}[T^+] = \frac{1}{2}(1 + 2 + \dots + n) = \frac{n(n+1)}{4},$$

and likewise for T^- .

Where do the tables come from?

Under H_0 each of the n ranks is positive or negative independently with probability $\frac{1}{2}$, giving 2^n equally likely sign patterns. The critical value at level α is the largest t with $\mathbb{P}(T^+ \leq t) \leq \alpha$.

Exercise. For $n = 4$, list all $2^4 = 16$ sign patterns and find the distribution of T^+ . Show that *no* significant result is possible at the 5% level, so the booklet table has no row for $n = 4$. Then verify the first row of the table ($n = 5$).

With ranks 1, 2, 3, 4 ($T^+ + T^- = 10$), counting which subsets of $\{1, 2, 3, 4\}$ give each sum:

t	0	1	2	3	4	5	6	7	8	9	10
$16\mathbb{P}(T^+ = t)$	1	1	1	2	2	2	2	2	1	1	1

(e.g. $T^+ = 3$ from $\{3\}$ or $\{1, 2\}$.) The most extreme possible result is $T = 0$, with

$$\mathbb{P}(T^+ \leq 0) = \frac{1}{16} = 0.0625 > 0.05,$$

so even a perfect run of signs is not significant at 5%: no test is possible with $n = 4$.

For $n = 5$ there are $2^5 = 32$ patterns and

$$\mathbb{P}(T^+ \leq 0) = \frac{1}{32} = 0.03125 \leq 0.05, \quad \mathbb{P}(T^+ \leq 1) = \frac{2}{32} = 0.0625 > 0.05,$$

so the critical value at the one-tailed 5% level is $T = 0$ — and since $0.03125 > 0.025$ there is no entry at the 2.5% level. This is exactly the first row of the formula-booklet table.

Example (OCR Further Statistics, December 2018)

The reaction times, in milliseconds, of all adult males in a standard experiment have a symmetrical distribution with mean and median both equal to 700 and standard deviation 125. The reaction times of a random sample of 6 international athletes are measured and the results are as follows:

702, 631, 540, 714, 575, 480.

It is required to test whether international athletes have a mean reaction time which is less than 700.

- (a) Assume first that the reaction times of international athletes have the distribution $N(\mu, 125^2)$. Test at the 5% significance level whether $\mu < 700$.
- (b) Now assume only that the distribution of the data is symmetrical, but not necessarily normal.
 - (i) State with a reason why a Wilcoxon test is preferable to a sign test.
 - (ii) Use an appropriate Wilcoxon test at the 5% significance level to test whether the median reaction time of international athletes is less than 700.
- (c) Explain why the significance tests in part (a) and part (b)(ii) could produce different results.

- (a) $H_0: \mu = 700$; $H_1: \mu < 700$, where μ is the mean reaction time of all international athletes. The sample mean is $\bar{x} = \frac{702+631+540+714+575+480}{6} = 607$, so

$$z = \frac{607 - 700}{125/\sqrt{6}} = -1.822, \quad p = \mathbb{P}(Z < -1.822) = 0.0342.$$

Since $0.0342 < 0.05$ (equivalently $-1.822 < -1.645$), we reject H_0 : there is evidence at the 5% level to suggest that the mean reaction time of international athletes is less than 700 ms.

- (b) (i) The Wilcoxon test uses the magnitudes of the differences from 700 as well as their signs, so it uses more of the information in the data.
(ii) $H_0: m = 700$; $H_1: m < 700$, where m is the median reaction time of all international athletes. The differences from 700 are

$$+2, -69, -160, +14, -125, -220,$$

with ranks of $|d_i|$ equal to 1, 3, 5, 2, 4, 6 respectively. So

$$T^+ = 1 + 2 = 3, \quad T^- = 3 + 5 + 4 + 6 = 18 \quad (3 + 18 = 21 = \frac{6 \times 7}{2} \checkmark)$$

and the test statistic is $T = T^+ = 3$. From the booklet, the critical value for $n = 6$ at the one-tailed 5% level is 2. Since $T = 3 > 2$, we do not reject H_0 : there is insufficient evidence at the 5% level to suggest that the median reaction time of international athletes is less than 700 ms.

- (c) The two tests make different assumptions about the population (normal with $\sigma = 125$ versus symmetric only) and use the data differently, so they can reach different conclusions — as they do here: the z-test extracts more from the data, at the price of a stronger assumption.

Textbook Exercises: [CUPS] Ch 4 §2; [S3&4] S4 Ch 2

Paired Samples

Often two measurements are made on the *same* subjects — before and after training, two treatments on matched patients, two judges scoring the same performances. Such data is **paired**. The right move is to work with the *differences* within each pair, reducing the problem to a single sample of differences.

Tip (Paired or two-sample?)

If each value in one sample is naturally linked to one particular value in the other (same person, same plot of land, same day), use a **paired** test on the differences. If the two samples are independent groups (possibly of different sizes), a paired test is impossible — use the two-sample rank-sum test of the next section.

Fact (The correct null hypothesis) — For a paired test the null hypothesis is

$$H_0: m_d = 0,$$

where m_d is the **population median of the differences** in the given context. It is **not** “the two medians are equal”: the median of the differences is not the difference of the medians.

Example

Find paired data sets X and Y with equal medians but $m_d \neq 0$.

Take the pairs $(x, y) = (2, 5), (5, 1), (8, 5)$. Then X -values 2, 5, 8 have median 5 and Y -values 5, 1, 5 have median 5 — equal medians. But the differences $x - y$ are $-3, 4, 3$, with median 3 $\neq 0$. (Although seen in some older exam papers, $H_0: m_X = m_Y$ is simply wrong.)

Example (Paired sign test)

Ten athletes run a time trial before and after a training programme. The differences (before – after, seconds) are

$$+1.2, +0.8, -0.3, +2.1, +0.5, 0, +1.7, +0.9, +0.4, +1.1.$$

Use a sign test at the 5% significance level to decide whether the programme reduces times.

$H_0: m_d = 0$; $H_1: m_d > 0$, where m_d is the population median of the differences (before – after).

Discard the zero difference, leaving $n = 9$ with 8 positive signs. Under H_0 , $S \sim B(9, \frac{1}{2})$:

$$\mathbb{P}(S \geq 8) = \frac{\binom{9}{8} + \binom{9}{9}}{2^9} = \frac{9 + 1}{512} = 0.0195$$

Since $0.0195 < 0.05$, we reject H_0 . There is evidence at the 5% level to suggest that the training programme reduces the athletes' times.

Example (Paired Wilcoxon signed-rank test)

Eight students sit a paper before and after a revision course. The score differences (after – before) are

$$+3, +7, -2, +9, +5, +11, -1, +6.$$

Assuming the differences are symmetrically distributed, test at the 5% significance level whether the course improves scores.

$H_0: m_d = 0; H_1: m_d > 0$, where m_d is the population median score difference (after – before).

Ranking the $|d_i|$: the values 1, 2, 3, 5, 6, 7, 9, 11 receive ranks 1, 2, 3, 4, 5, 6, 7, 8 respectively. The negative differences are -1 (rank 1) and -2 (rank 2), so

$$T^- = 1 + 2 = 3, \quad T^+ = 36 - 3 = 33.$$

Under H_1 we expect T^- to be small, so the test statistic is $T = T^- = 3$. From the booklet, the critical value for $n = 8$ at the one-tailed 5% level is 5.

Since $T = 3 \leq 5$, we reject H_0 . There is evidence at the 5% level to suggest that the revision course improves scores.

Example (OCR S4, June 2014)

A teacher believes that the calculator paper in a GCSE Mathematics examination was easier than the non-calculator paper. The marks of a random sample of ten students are shown in the table.

Student	A	B	C	D	E	F	G	H	I	J
Mark on paper 1 (non-calculator)	66	79	58	87	67	55	75	62	50	84
Mark on paper 2 (calculator)	57	84	70	90	75	42	82	72	65	82

- (i) Use a Wilcoxon signed-rank test, at the 5% significance level, to test the teacher’s belief.
- (ii) State the assumption necessary for this test to be applied.

(i) $H_0: m_d = 0; H_1: m_d > 0$, where m_d is the population median of the differences (paper 2 – paper 1): “easier” means higher marks on the calculator paper.

Student	A	B	C	D	E	F	G	H	I	J
d_i (paper 2 – paper 1)	-9	+5	+12	+3	+8	-13	+7	+10	+15	-2
rank of $ d_i $	6	3	8	2	5	9	4	7	10	1

$$T^- = 6 + 9 + 1 = 16, \quad T^+ = 55 - 16 = 39.$$

Under H_1 we expect T^- to be small, so the test statistic is $T = T^- = 16$. From the booklet, the critical value for $n = 10$ at the one-tailed 5% level is 10.

Since $T = 16 > 10$, we do not reject H_0 . There is insufficient evidence at the 5% level to suggest that the calculator paper was easier than the non-calculator paper.

- (ii) The population of differences is symmetrically distributed.

Textbook Exercises: [CUP.S] Ch 4 §3; [S3&4] S4 Ch 2

The Wilcoxon Rank-Sum Test

For two *independent* (unpaired) samples, of sizes m and n with $m \leq n$, we use the **Wilcoxon rank-sum test**, also known as the **Mann–Whitney U test**. The hypotheses are

H_0 : the two population distributions are identical

H_1 : the two population distributions are not identical.

More commonly we assume the two distributions have the same *shape* and may differ only in *location*, so the hypotheses become $H_0: m_A = m_B$ against $H_1: m_A \neq m_B$ (or one-tailed), where m_A, m_B are the population medians.

- Tip (The procedure)**
1. Combine both samples into a single list and rank all $m + n$ values from 1 (smallest) upwards.
 2. Let W be the sum of the ranks of the **smaller** sample (size m).
 3. Reject H_0 if $W \leq W_{\text{crit}}$ (booklet table) or $W \geq m(m + n + 1) - W_{\text{crit}}$. For a one-tailed test, use whichever tail H_1 points to.

Remark. The Mann–Whitney statistic is $U = W - \frac{m(m+1)}{2}$, the number of (smaller-sample, larger-sample) pairs in which the smaller-sample value wins; some tables are written in terms of U instead of W . OCR uses W .

Example (Rank-sum test)

Two groups of students learn a routine by different methods, then are timed performing it. The times (minutes) are

Method A ($m = 4$): 11, 12, 15, 19

Method B ($n = 6$): 17, 18, 21, 24, 25, 30

Test at the 5% significance level whether the median times under the two methods differ. (You may assume the two distributions have the same shape.)

$H_0: m_A = m_B$; $H_1: m_A \neq m_B$, where m_A, m_B are the population median times for the two methods.
Ranking all 10 values together:

value	11	12	15	17	18	19	21	24	25	30
rank	1	2	3	4	5	6	7	8	9	10
sample	A	A	A	B	B	A	B	B	B	B

The smaller sample is A ($m = 4, n = 6$):

$$W = 1 + 2 + 3 + 6 = 12.$$

From the booklet, the lower critical value for $m = 4, n = 6$ at the 5% two-tailed level is 12 (and the upper is $m(m + n + 1) - 12 = 44 - 12 = 32$).

Since $W = 12 \leq 12$, we reject H_0 . There is evidence at the 5% level to suggest that the median times under the two methods differ.

Fact — Under H_0 , the distribution of W is symmetric about its mean

$$\mathbb{E}[W] = \frac{m(m+n+1)}{2}.$$

The proof is a pleasant symmetry argument.

The smallest W can be is $1 + 2 + \dots + m = \frac{m(m+1)}{2}$ (the smaller sample takes all the lowest ranks); the largest is $(n+1) + (n+2) + \dots + (n+m) = \frac{m(m+1)}{2} + mn$. Re-ranking the data in the opposite order replaces each rank r by $m+n+1-r$, hence replaces W by $m(m+n+1) - W$; under H_0 both rankings are equally valid, so the distribution of W is symmetric about

$$\mathbb{E}[W] = \frac{1}{2} \left[\frac{m(m+1)}{2} + \frac{m(m+1)}{2} + mn \right] = \frac{m(m+n+1)}{2}.$$

Where do the tables come from?

Example ($m = 2, n = 5$)

Under H_0 all $\binom{7}{2} = 21$ choices of the two ranks held by the smaller sample are equally likely. Find the distribution of W and hence the one-tailed critical value at the 5% level.

The possible values of W run from $1 + 2 = 3$ up to $6 + 7 = 13$. Counting pairs of distinct ranks from $\{1, \dots, 7\}$ with each sum:

w	3	4	5	6	7	8	9	10	11	12	13
$21 \mathbb{P}(W = w)$	1	1	2	2	3	3	3	2	2	1	1

Note the symmetry about $\mathbb{E}[W] = \frac{2 \times 8}{2} = 8$, as predicted. Now

$$\mathbb{P}(W \leq 3) = \frac{1}{21} = 0.048 \leq 0.05, \quad \mathbb{P}(W \leq 4) = \frac{2}{21} = 0.095 > 0.05,$$

so the one-tailed 5% critical value is 3: only the most extreme arrangement possible is significant. Every entry in the booklet table is computed by exactly this kind of enumeration.

Textbook Exercises: [CUP.S] Ch 4 §4; [S3&4] S4 Ch 2

Normal Approximations for Large Samples

The booklet tables stop at modest sample sizes. For larger samples, T^+ and W are sums of many small independent contributions, so (by CLT-style reasoning) they are approximately normal, with the means we have already found.

Fact (In the formula booklet) — For large samples, under H_0 :

$$\text{Wilcoxon signed-rank: } T \sim N\left(\frac{n(n+1)}{4}, \frac{n(n+1)(2n+1)}{24}\right) \text{ approximately}$$

$$\text{Wilcoxon rank-sum: } W \sim N\left(\frac{m(m+n+1)}{2}, \frac{mn(m+n+1)}{12}\right) \text{ approximately}$$

Since T and W are discrete (integer-valued) and the normal distribution is continuous, a **continuity correction** must be used. A continuity correction is essential here.

The means are the ones we derived earlier; the variance of T^+ is no mystery either.

Rank i contributes i to T^+ with probability $\frac{1}{2}$ and 0 otherwise, independently of the other ranks, so its contribution X_i has

$$\text{Var}[X_i] = \mathbb{E}[X_i^2] - \mathbb{E}[X_i]^2 = \frac{i^2}{2} - \frac{i^2}{4} = \frac{i^2}{4}.$$

Summing these independent contributions, using $\sum_{i=1}^n i^2 = \frac{n(n+1)(2n+1)}{6}$:

$$\text{Var}[T^+] = \frac{1}{4} \sum_{i=1}^n i^2 = \frac{n(n+1)(2n+1)}{24}.$$

Remark. The approximations are for T^+ (equivalently T^-) and W , so we can test at either tail, choosing the tail H_1 points to; for a two-tailed test, double the tail probability.

Example (Signed-rank, large sample)

A website claims the median time to complete its checkout is 35 seconds. The times of 20 customers are recorded; none equals 35, and the Wilcoxon signed-rank calculation gives $T^+ = 60$ (sum of ranks of times above 35). Test at the 5% significance level whether the median time is less than claimed.

$H_0: m = 35; \quad H_1: m < 35$, where m is the population median checkout time.

With $n = 20$, under H_0 :

$$\mathbb{E}[T^+] = \frac{20 \times 21}{4} = 105, \quad \text{Var}[T^+] = \frac{20 \times 21 \times 41}{24} = 717.5,$$

so $T^+ \approx N(105, 717.5)$. Under H_1 we expect T^+ small. With a continuity correction ($T^+ \leq 60$ becomes $T^+ < 60.5$):

$$\mathbb{P}(T^+ \leq 60) \approx \mathbb{P}\left(Z < \frac{60.5 - 105}{\sqrt{717.5}}\right) = \mathbb{P}(Z < -1.661) = 0.0483$$

Since $0.0483 < 0.05$, we reject H_0 . There is evidence at the 5% level to suggest that the median checkout time is less than 35 seconds.

Example (OCR S4, June 2018)

A Wilcoxon signed-rank test is carried out at the 5% level of significance on a random sample of size 32. The hypotheses are $H_0: m = m_0$, $H_1: m < m_0$, where m is the population median and m_0 is a specific numerical value. The value obtained for the test statistic T is 162. Find the outcome of the test.

With $n = 32$, under H_0 :

$$\mathbb{E}[T] = \frac{32 \times 33}{4} = 264, \quad \text{Var}[T] = \frac{32 \times 33 \times 65}{24} = 2860,$$

so $T \approx N(264, 2860)$. Under H_1 we expect T to be small, and $162 < 264$. With a continuity correction:

$$\mathbb{P}(T \leq 162) \approx \mathbb{P}\left(Z < \frac{162.5 - 264}{\sqrt{2860}}\right) = \mathbb{P}(Z < -1.898) = 0.0289$$

Since $0.0289 < 0.05$, we reject H_0 . There is evidence at the 5% level to suggest that the population median is less than m_0 .

Example (Rank-sum, large sample)

Independent samples of $m = 10$ and $n = 12$ plants are grown with two fertilisers and their heights ranked together. The rank sum of the smaller sample is $W = 85$. Test at the 5% significance level whether the two fertilisers produce different median heights.

$H_0: m_A = m_B$; $H_1: m_A \neq m_B$ (medians of the two height distributions).

Under H_0 :

$$\mathbb{E}[W] = \frac{10 \times 23}{2} = 115, \quad \text{Var}[W] = \frac{10 \times 12 \times 23}{12} = 230,$$

so $W \approx N(115, 230)$. Observed $W = 85$ is below the mean; with a continuity correction:

$$\mathbb{P}(W \leq 85) \approx \mathbb{P}\left(Z < \frac{85.5 - 115}{\sqrt{230}}\right) = \mathbb{P}(Z < -1.945) = 0.0259$$

The test is two-tailed, so the p -value is $2 \times 0.0259 = 0.0518 > 0.05$.

We do not reject H_0 : there is insufficient evidence at the 5% level to suggest that the two fertilisers produce different median heights. (Had the test been one-tailed, the conclusion would have been different — read the question carefully.)

Choosing the right test

Situation	Test	Assumptions
One sample, median	Sign test	none
One sample, median	Wilcoxon signed-rank	symmetric distribution
Paired samples	Sign test on differences	none
Paired samples	Wilcoxon signed-rank on differences	symmetric differences
Two independent samples	Wilcoxon rank-sum	same shape distributions

Tip

In conclusions, never write “accept H_0 ” or “accept H_1 ”. Either “reject H_0 : there is evidence at the [level] to suggest that...” or “do not reject H_0 : there is insufficient evidence to suggest that...” — always in the context of the question.

Textbook Exercises: [CUPS] Ch 4 §5; [S3&4] S4 Ch 2