

## Chapter 511

# Non-Inferiority Tests for Pairwise Mean Differences in a Williams Cross-Over Design

---

## Introduction

Senn (2002) defines a *cross-over* design as one in which each subject receives all treatments, and the objective is to study differences among the treatments. The name *cross-over* comes from the most common case in which there are only two treatments. In this case, each subject *crosses over* from one treatment to the other. It is assumed that there is a *washout* period between treatments during which the response returns back to its baseline value. If this does not occur, there is said to be a *carry-over* effect.

An  $a \times k$  cross-over design contains  $a$  sequences (treatment orderings) and  $k$  time periods (occasions) corresponding to the  $k$  treatments. The design includes a washout period between responses to make certain that the effects of the first drug do not carry over to the second. Thus, the groups in this design are defined by the sequence in which the drugs are administered, not by the treatments they receive.

Cross-over designs are employed because, if the no-carryover assumption is met, treatment differences are measured within a subject rather than between subjects—making a more precise measurement. Examples of the situations that might use a cross-over design are the comparison of anti-inflammatory drugs in arthritis and the comparison of hypotensive agents in essential hypertension. In both cases, symptoms are expected to return to their usual baseline level shortly after the treatment is stopped.

The sample size calculations in the procedure are based on the formulas presented in Chow, Shao, Wang, & Lohknygina (2018).

---

## Advantages of Cross-Over Designs

A comparison of treatments on the same subject is expected to be more precise. The increased precision often translates into a smaller sample size. Also, patient enrollment into the study may be easier because each patient will receive both treatments. Finally, it is often more difficult to obtain a subject than to obtain a measurement.

---

## Disadvantages of Cross-Over Designs

The statistical analysis of a cross-over experiment is more complex than a parallel-group experiment and requires additional assumptions. It may be difficult to separate the treatment effect from the period effect, the carry-over effect of the previous treatment, and the interaction between period and treatment.

The design cannot be used when the treatment (or the measurement of the response) alters the subject permanently. Hence, it should not be used to compare treatments that are intended to provide a cure.

Because subjects must be measured at least twice, it is often more difficult to keep patients enrolled in the study. It is arguably simpler to measure a subject once than to obtain their measurement twice. This is particularly true when the measurement process is painful, uncomfortable, embarrassing, or time consuming.

## Technical Details

The  $a \times k$  crossover design may be described as follows. Randomly assign the subjects to one of  $a$  sequence groups with  $n_1$  subjects in sequence one,  $n_2$  subjects in sequence two, and so forth up to sequence  $a$ . In order to achieve design balance, the sample sizes  $n_1, n_2, \dots, n_a$  are assumed to be equal so that  $n_1 = n_2 = \dots = n_a = n = N/a$ . Sequence one is given a specific sequence of  $k$  treatments, sequence two is given a different sequence of the same  $k$  treatments, and so forth up to sequence  $a$ .

### Williams Cross-Over Design

Williams cross-over designs are constructed from Latin squares as outlined in Chow and Liu (2009). If the number of treatments ( $k$ ) is even, then Williams design results in a  $k \times k$  cross-over design (i.e., with  $k$  sequences and  $k$  treatments/periods). If the number of treatments ( $k$ ) is odd, then Williams design results in a  $2k \times k$  cross-over design (i.e., with  $2k$  sequences and  $k$  treatments/periods). For example, a Williams design with 4 treatments would result in a  $4 \times 4$  cross-over design and would have 4 sequences with 4 periods corresponding to the 4 treatments. On the other hand, a Williams design with 3 treatments would result in a  $6 \times 3$  cross-over design and would have 6 sequences with 3 periods corresponding to the 3 treatments.

Define  $y_{ijl}$  as the continuous response from subject  $j$  ( $j = 1, \dots, n$ ) in sequence  $i$  ( $i = 1, \dots, a$ ) given treatment  $l$  ( $l = 1, \dots, k$ ). The observations taken from the same subject may be correlated with one another.

Further define the paired differences between treatments  $u$  and  $v$  for each subject within each sequence as

$$d_{ij}(u, v) = y_{iju} - y_{ijv}$$

and the overall true difference as

$$\delta = \mu_u - \mu_v.$$

The overall difference can be estimated as

$$\hat{\delta} = \frac{1}{an} \sum_{i=1}^a \sum_{j=1}^n d_{ij}(u, v).$$

The estimated difference is asymptotically normally distributed with variance  $\sigma_d^2$ , which can be estimated as

$$\hat{\sigma}_d^2 = \frac{1}{a(n-1)} \sum_{i=1}^a \sum_{j=1}^n (d_{ij}(u, v) - \bar{d}_{i\cdot}(u, v))^2,$$

where

$$\bar{d}_{i\cdot}(u, v) = \frac{1}{n} \sum_{j=1}^n d_{ij}(u, v).$$

## Non-Inferiority Tests for Pairwise Mean Differences in a Williams Cross-Over Design

The standard deviation, then, is

$$SD = \sigma_d = \sqrt{\sigma_d^2}$$

with estimate

$$\widehat{SD} = \hat{\sigma}_d = \sqrt{\hat{\sigma}_d^2}.$$

---

## Non-Inferiority Test Statistics

### Higher Means Better

When higher means are better, the null and alternative hypotheses for a one-sided non-inferiority test are

$$H_0: \mu_u - \mu_v \leq D_0 \quad \text{vs.} \quad H_A: \mu_u - \mu_v > D_0$$

or equivalently

$$H_0: \delta \leq D_0 \quad \text{vs.} \quad H_A: \delta > D_0$$

where  $D_0$  is the lower non-inferiority bound (i.e., the smallest difference  $(\mu_u - \mu_v)$  for which treatment  $u$  will still be considered non-inferior to treatment  $v$ ). When higher means are better,  $D_0$  should be less than zero.

The power and sample size calculations are based on the test statistic

$$t = \frac{\hat{\delta} - D_0}{\frac{\hat{\sigma}_d}{\sqrt{an}}}$$

which follows a central  $T$  distribution with  $a(n - 1)$  degrees of freedom under the null hypothesis. The null hypothesis is rejected in favor of the alternative at level  $\alpha$  if

$$\frac{\hat{\delta} - D_0}{\frac{\hat{\sigma}_d}{\sqrt{an}}} > t_{1-\alpha, a(n-1)}$$

where  $t_{1-\alpha, a(n-1)}$  is the upper  $1 - \alpha$  percentile of a central  $T$  distribution with  $a(n - 1)$  degrees of freedom.

## Higher Means Worse

When higher means are worse, the null and alternative hypotheses for a one-sided non-inferiority test are

$$H_0: \mu_u - \mu_v \geq D_0 \quad \text{vs.} \quad H_A: \mu_u - \mu_v < D_0$$

or equivalently

$$H_0: \delta \geq D_0 \quad \text{vs.} \quad H_A: \delta < D_0$$

where  $D_0$  is the upper non-inferiority bound (i.e., the largest difference  $(\mu_u - \mu_v)$  for which treatment  $u$  will still be considered non-inferior to treatment  $v$ ). When higher means are worse,  $D_0$  should be greater than zero.

The power and sample size calculations are based on the test statistic

$$t = \frac{\hat{\delta} - D_0}{\frac{\hat{\sigma}_d}{\sqrt{an}}}$$

which follows a central  $T$  distribution with  $a(n - 1)$  degrees of freedom under the null hypothesis. The null hypothesis is rejected in favor of the alternative at level  $\alpha$  if

$$\frac{\hat{\delta} - D_0}{\frac{\hat{\sigma}_d}{\sqrt{an}}} < t_{\alpha, a(n-1)}$$

where  $t_{\alpha, a(n-1)}$  is the lower  $\alpha$  percentile of a central  $T$  distribution with  $a(n - 1)$  degrees of freedom.

## Bonferroni Adjustment for Multiple Tests

In a design with  $k$  treatments, there are  $k(k - 1)/2$  possible pairwise  $(u, v)$  comparison tests. To protect the overall alpha level, the individual test alpha level is often divided by the number of tests performed. This is known as the Bonferroni adjustment for multiple comparisons. When this adjustment is used in hypothesis testing, the individual test alpha value of  $\alpha/(k(k - 1)/2)$  is substituted for  $\alpha$  in the formulas above.

## Non-Inferiority Power Calculations

### Higher Means Better

According to Chow, Shao, Wang, & Lokhnygina (2018) page 65, the power for the one-sided non-inferiority test of  $H_0: \delta \leq D_0$  versus  $H_A: \delta > D_0$  is

$$1 - T_{a(n-1)} \left( t_{1-\alpha, a(n-1)} \left| \frac{\delta_1 - D_0}{\frac{\sigma_d}{\sqrt{an}}} \right| \right)$$

where  $T_{df}(X|NCP)$  is the non-central  $T$  distribution function with  $df$  degrees of freedom and non-centrality parameter  $NCP$  evaluated at  $X$ ,  $\delta_1$  is the actual value of the minimum difference under the alternative hypothesis, and  $t_{1-\alpha, a(n-1)}$  is the upper  $1 - \alpha$  percentile of a central  $T$  distribution with  $a(n - 1)$  degrees of freedom. The sample size is determined using a binary search of possible values for  $n$ .

### Higher Means Worse

Derived from Chow, Shao, Wang, & Lokhnygina (2018) page 65, the power for the one-sided non-inferiority test of  $H_0: \delta \geq D_0$  versus  $H_A: \delta < D_0$  is

$$1 - T_{a(n-1)} \left( t_{1-\alpha, a(n-1)} \left| \frac{D_0 - \delta_1}{\frac{\sigma_d}{\sqrt{an}}} \right| \right)$$

where  $T_{df}(X|NCP)$  is the non-central  $T$  distribution function with  $df$  degrees of freedom and non-centrality parameter  $NCP$  evaluated at  $X$ ,  $\delta_1$  is the actual value of the minimum difference under the alternative hypothesis, and  $t_{1-\alpha, a(n-1)}$  is the upper  $1 - \alpha$  percentile of a central  $T$  distribution with  $a(n - 1)$  degrees of freedom. The sample size is determined using a binary search of possible values for  $n$ .

### Bonferroni Adjustment for Multiple Tests

In a design with  $k$  treatments, there are  $k(k - 1)/2$  possible pairwise  $(u, v)$  comparison tests. To protect the overall alpha level, the individual test alpha level is often divided by the number of tests performed. This is known as the Bonferroni adjustment for multiple comparisons. When this adjustment is used in power calculations, the individual test alpha value of  $\alpha/(k(k - 1)/2)$  is substituted for  $\alpha$  in the formulas above.

# Example 1 – Power Analysis

Suppose you want to consider the power of a balanced Williams cross-over design with 3 groups and a continuous endpoint where the test is computed based on the difference for sequence sample sizes between 30 and 100. The actual minimum difference is 0, the non-inferiority difference is -0.5, and the estimated standard deviation of the paired differences is 3.5. The overall significance level is 0.05 with individual test alpha adjusted for 3 tests.

## Setup

If the procedure window is not already open, use the PASS Home window to open it. The parameters for this example are listed below and are stored in the **Example 1** settings file. To load these settings to the procedure window, click **Open Example Settings File** in the Help Center or File menu.

Design Tab	
Solve For .....	<b>Power</b>
Higher Means Are .....	<b>Better</b>
Alpha.....	<b>0.05</b>
Adjust Alpha for Multiple Tests .....	<b>Checked</b>
k (Number of Treatments).....	<b>3</b>
n (Sample Size per Sequence) .....	<b>30 to 100 by 10</b>
D0 (Non-Inferiority Difference) .....	<b>-0.5</b>
D1 (Minimum Difference H1) .....	<b>0</b>
Standard Deviation (SD).....	<b>3.5</b>

## Non-Inferiority Tests for Pairwise Mean Differences in a Williams Cross-Over Design

## Output

Click the Calculate button to perform the calculations and generate the following output.

## Numeric Results

Solve For: [Power](#)  
 Design: 6x3 Williams Cross-Over Design  
 Higher Proportions Are: Better  
 Hypotheses:  $H_0: \mu_u - \mu_v \leq D_0$  vs.  $H_1: \mu_u - \mu_v > D_0$  for  $u, v = 1, \dots, 3$  with  $u \neq v$ .  
 Number of Possible Tests: 3

Power	Sample Size		Difference		Standard Deviation SD	Alpha*	
	Sequence n	Total N	Non-Inferiority D0	Minimum D1		Overall	Individual Test
0.41142	30	180	-0.5	0	3.5	0.05	0.017
0.52964	40	240	-0.5	0	3.5	0.05	0.017
0.63186	50	300	-0.5	0	3.5	0.05	0.017
0.71695	60	360	-0.5	0	3.5	0.05	0.017
0.78572	70	420	-0.5	0	3.5	0.05	0.017
0.83997	80	480	-0.5	0	3.5	0.05	0.017
0.88191	90	540	-0.5	0	3.5	0.05	0.017
0.91380	100	600	-0.5	0	3.5	0.05	0.017

\* Alpha was adjusted for 3 tests using the Bonferroni method. Power was calculated using Individual Test Alpha.

Power The probability of rejecting a false null hypothesis when the alternative hypothesis is true.  
 n The sample size in each sequence.  
 N The total sample size from all 6 sequences combined. The sample is divided equally among sequences.  
 D0 The non-inferiority difference used to specify the hypothesis test.  
 D1 The minimum treatment difference to detect at which power is calculated.  $D1 = \text{Minimum of } (\mu_u - \mu_v) | H_1 \text{ for } u, v = 1, \dots, k \text{ with } u \neq v$ .  
 SD The standard deviation of paired differences. This is estimated from a previous study.  
 Alpha The probability of rejecting a true null hypothesis.

## Summary Statements

A 6x3 Williams cross-over design (6 sequences, 3 treatments) (where higher means are considered to be better) will be used to test whether each treatment mean is non-inferior to the others, with a non-inferiority margin of -0.5 ( $H_0: \mu_u - \mu_v \leq -0.5$  versus  $H_1: \mu_u - \mu_v > -0.5$ , for  $u, v = 1, \dots, 3$  with  $u \neq v$ ). Each comparison will be made using a one-sided t-test. The Type I error rate ( $\alpha$ ) for each comparison is 0.017 (Bonferroni-adjusted for 3 comparisons), and the overall Type I error rate across all tests is 0.05. The standard deviation of paired differences is assumed to be 3.5. To detect a minimum mean difference ( $\mu_u - \mu_v$ ) of 0 with a sample size of 30 in each sequence (totaling 180 subjects), the power is 0.41142.

## Non-Inferiority Tests for Pairwise Mean Differences in a Williams Cross-Over Design

## Dropout-Inflated Sample Size

Group	Dropout Rate	Sample Size Ni	Dropout- Inflated Enrollment Sample Size Ni'	Expected Number of Dropouts Di
1 - 6	20%	30	38	8
Total		180	228	48
1 - 6	20%	40	50	10
Total		240	300	60
1 - 6	20%	50	63	13
Total		300	378	78
1 - 6	20%	60	75	15
Total		360	450	90
1 - 6	20%	70	88	18
Total		420	528	108
1 - 6	20%	80	100	20
Total		480	600	120
1 - 6	20%	90	113	23
Total		540	678	138
1 - 6	20%	100	125	25
Total		600	750	150

Group Lists the group numbers.

Dropout Rate The percentage of subjects (or items) that are expected to be lost at random during the course of the study and for whom no response data will be collected (i.e., will be treated as "missing"). Abbreviated as DR.

Ni The evaluable sample size for each group at which power is computed (as entered by the user). If Ni subjects are evaluated out of the Ni' subjects that are enrolled in the study, the design will achieve the stated power.

Ni' The number of subjects that should be enrolled in each group in order to obtain Ni evaluable subjects, based on the assumed dropout rate. Ni' is calculated by inflating Ni using the formula  $Ni' = Ni / (1 - DR)$ , with Ni' always rounded up. (See Julious, S.A. (2010) pages 52-53, or Chow, S.C., Shao, J., Wang, H., and Lohknygina, Y. (2018) pages 32-33.)

Di The expected number of dropouts in each group.  $Di = Ni' - Ni$ .

## Dropout Summary Statements

Anticipating a 20% dropout rate, group sizes of 38, 38, 38, 38, 38, and 38 subjects should be enrolled to obtain final group sample sizes of 30, 30, 30, 30, 30, and 30 subjects.

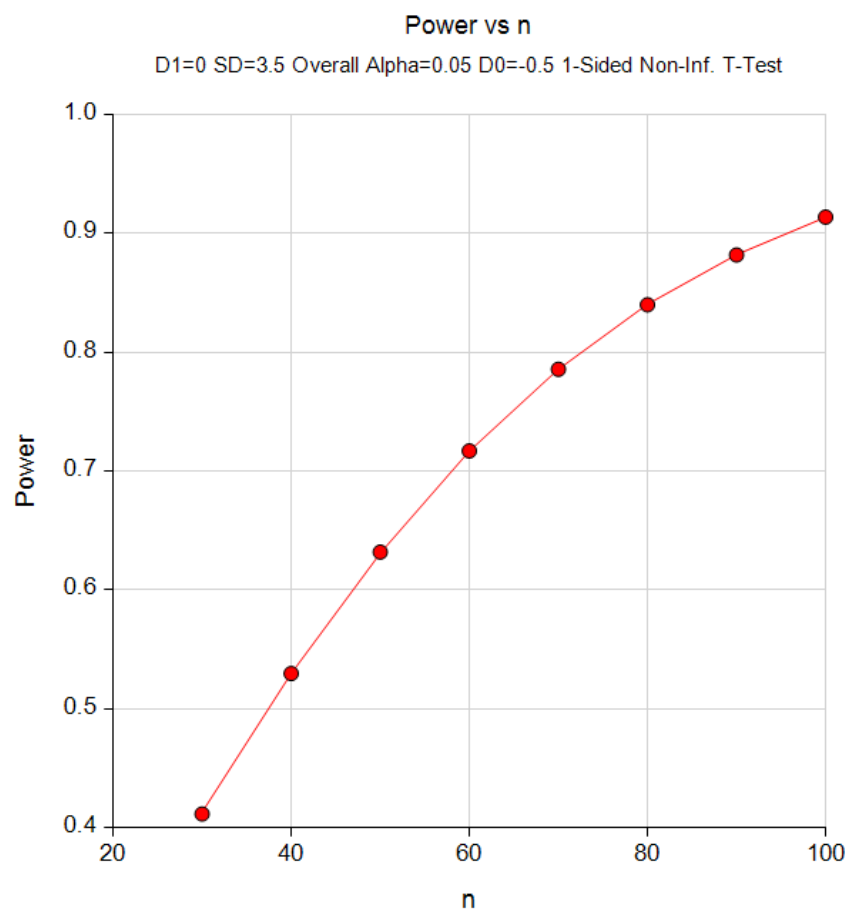
## References

Chow, S.C., Shao, J., Wang, H., and Lohknygina, Y. 2018. Sample Size Calculations in Clinical Research, Third Edition. Taylor & Francis/CRC. Boca Raton, Florida.



## Non-Inferiority Tests for Pairwise Mean Differences in a Williams Cross-Over Design

## Plots



This report shows the values of each of the parameters, one scenario per row. This plot shows the relationship between sample size and power. We see that a sample size of just over 70 per sequence is required to detect a minimum difference of 0 with 80% power when the lower non-inferiority bound is -0.5.

## Example 2 – Calculating Sample Size (Validation using Hand Calculations)

In this example, we'll find the sample size required in a  $6 \times 3$  Williams cross-over design ( $k = 3$ ) to detect a difference of -0.05 with 80% power in a non-inferiority test with a margin of -0.5 with a significance level of 0.05 when the standard deviation of paired differences is 1.5. We'll make no adjustment for multiple testing in this example. We'll also validate this procedure by computing power values manually.

The power for per-sequence sample sizes of 11 and 12 calculated by hand using the power formula referenced earlier is

$$\text{Power} = 1 - T_{a(n-1)} \left( t_{1-\alpha, a(n-1)} \left| \frac{\delta_1 - D_0}{\frac{\sigma_d}{\sqrt{an}}} \right| \right)$$

$$\begin{aligned} \text{Power}_{(n=11)} &= 1 - T_{60} \left( 1.670649 \left| \frac{-0.05 + 0.5}{\frac{1.5}{\sqrt{6 \times 11}}} \right| \right) \\ &= 0.777782 \end{aligned}$$

$$\begin{aligned} \text{Power}_{(n=12)} &= 1 - T_{66} \left( 1.668271 \left| \frac{-0.05 + 0.5}{\frac{1.5}{\sqrt{6 \times 12}}} \right| \right) \\ &= 0.809076 \end{aligned}$$

These results indicate that the minimum required sample size per group is 12, since it is the smallest sample size that achieves the desired 80% power.

## Setup

If the procedure window is not already open, use the PASS Home window to open it. The parameters for this example are listed below and are stored in the **Example 2** settings file. To load these settings to the procedure window, click **Open Example Settings File** in the Help Center or File menu.

### Design Tab

Solve For ..... **Sample Size**  
 Higher Means Are ..... **Better**  
 Power ..... **0.80**  
 Alpha ..... **0.05**  
 Adjust Alpha for Multiple Tests ..... **Unchecked**  
 k (Number of Treatments) ..... **3**  
 D0 (Non-Inferiority Difference) ..... **-0.5**  
 D1 (Minimum Difference|H1) ..... **-0.05**  
 Standard Deviation (SD) ..... **1.5**

## Output

Click the Calculate button to perform the calculations and generate the following output.

### Numeric Results

Solve For: [Sample Size](#)  
 Design: 6x3 Williams Cross-Over Design  
 Higher Proportions Are: Better  
 Hypotheses:  $H_0: \mu_u - \mu_v \leq D_0$  vs.  $H_1: \mu_u - \mu_v > D_0$  for  $u, v = 1, \dots, 3$  with  $u \neq v$ .  
 Number of Possible Tests: 3

Power	Sample Size		Difference		Standard Deviation SD	Alpha*
	Sequence n	Total N	Non-Inferiority D0	Minimum D1		
0.80908	12	72	-0.5	-0.05	1.5	0.05

\* Alpha was not adjusted for multiple tests.

The result from **PASS** match our hand calculations exactly.