Chapter 175

# Non-Inferiority Tests for the Difference of Two Proportions in a 2x2 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.

A 2×2 cross-over design contains two *sequences* (treatment orderings) and two time periods (occasions). One sequence receives treatment A followed by treatment B. The other sequence receives B and then A. 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. Indeed, higher-order cross-over designs have been used in which the same treatment is used on both occasions.

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, & Lokhnygina (2018) and Lui (2016). These sources discuss sample size calculations for 2×2 cross-over designs and provide power calculations methods that differ only in the way that the standard deviation is estimated.

## 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 2×2 crossover design may be described as follows. Randomly assign the subjects to one of two sequence groups so that there are $n_1$ subjects in sequence one and $n_2$ subjects in sequence two. In order to achieve design balance, the sample sizes $n_1$ and $n_2$ are assumed to be equal so that $n_1 = n_2 = N/2$.

Sequence one is given treatment A followed by treatment B. Sequence two is given treatment B followed by treatment A.

The design can be analyzed using a simple $z$-test if we ignore period and sequence effects or using a more complex random effects logistic regression model that adjusts for period and sequence effects. The sample size calculations herein ignore period and sequence effects. Julious (2010) suggests on page 175 that the bias due to ignoring period effects if a period-adjusted analysis is planned is not great and that sample size calculations that ignore period effects are adequate.

## Cross-Over Design

The discussions that follow summarize the results in Chow, Shao, Wang, & Lokhnygina (2018). Consider a 2×2 cross-over design and let $x_{ijk}$ represent the binary response (0 or 1) from the $j^{th}$ subject ($j$ = 1, ..., $n_i$) in the $i^{th}$ sequence ($i$ = 1, 2) given the $k^{th}$ treatment ($k$ = 1, 2). Here we assume that the sample sizes are equal in both sequences such that $n_1 = n_2 = n$. If replicates are taken from each subject (as in a 2×2$m$ replicated cross-over design) then $x_{ijk} = \bar{x}_{ijk.} = \frac{1}{m}\left(\sum_{l=1}^{m} x_{ijkl}\right)$, where $x_{ijkl}$ represents the $l^{th}$ binary response replicate ($l$ = 1, ..., $m$) from the $j^{th}$ subject ($j$ = 1, ..., $n_i$) in the $i^{th}$ sequence ($i$ = 1, 2) given the $k^{th}$ treatment ($k$ = 1, 2). The observations taken from the same subject may be correlated with one another. If we assume no sequence and period effects, then we can state that $P(x_{ijk} = 1) = P_k$. Further define the paired differences of treatment – control for each subject within each sequence as

$$d_{ij} = x_{ijT} - x_{ijC}$$

$$= x_{ij1} - x_{ij2}$$

and the overall treatment – control difference as

$$\delta = P_T - P_C$$

$$= P_1 - P_2.$$

The overall difference can be estimated as

$$\hat{\delta} = \frac{1}{2n} \sum_{i=1}^{2} \sum_{j=1}^{n} d_{ij}.$$

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

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

where

$$\bar{d}_{i.} = \frac{1}{n} \sum_{j=1}^{n} d_{ij}.$$

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 Proportions Better

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

$$H_0: P_T - P_C \leq D_0 \quad \text{vs.} \quad H_A: P_T - P_C > D_0$$

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 $(P_T - P_C)$ for which the treatment will still be considered non-inferior to the standard or control). When higher proportions are better, $D_0$ should be less than zero.

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

$$Z = \frac{\hat{\delta} - D_0}{\frac{\hat{\sigma}_d}{\sqrt{2n}}}$$

which is asymptotically distributed as standard normal under the null hypothesis. The null hypothesis is rejected in favor of the alternative at level $\alpha$ if

$$\frac{\hat{\delta} - D_0}{\frac{\hat{\sigma}_d}{\sqrt{2n}}} > Z_{1-\alpha}$$

where $Z_{1-\alpha}$ is the upper $1 - \alpha$ percentile of the standard normal distribution.

## Higher Proportions Worse

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

$$H_0: P_T - P_C \geq D_0 \quad \text{vs.} \quad H_A: P_T - P_C < D_0$$

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 $(P_T - P_C)$ for which the treatment will still be considered non-inferior to the standard or control). When higher proportions are worse, $D_0$ should be greater than zero.

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

$$Z = \frac{\hat{\delta} - D_0}{\frac{\hat{\sigma}_d}{\sqrt{2n}}}$$

which is asymptotically distributed as standard normal under the null hypothesis. The null hypothesis is rejected in favor of the alternative at level $\alpha$ if

$$\frac{\hat{\delta} - D_0}{\frac{\hat{\sigma}_d}{\sqrt{2n}}} < Z_\alpha$$

where $Z_\alpha$ is the lower $\alpha$ percentile of the standard normal distribution.

# Non-Inferiority Power Calculations

## Higher Proportions Better

According to Chow, Shao, Wang, & Lokhnygina (2018) page 84, the power for the one-sided non-inferiority test of $H_0: \delta \leq D_0$ versus $H_A: \delta > D_0$ is

$$\Phi\left(\frac{\delta_1 - D_0}{\frac{\sigma_d}{\sqrt{2n}}} - Z_{1-\alpha}\right)$$

where $\Phi()$ is the standard normal distribution function, $\delta_1$ is the actual value of the difference under the alternative hypothesis, and $Z_{1-\alpha}$ is the upper $1 - \alpha$ percentile of the standard normal distribution. The sample size is determined using a binary search of possible values for $n$.

## Higher Proportions Worse

Derived from Chow, Shao, Wang, & Lokhnygina (2018) page 84, the power for the one-sided non-inferiority test of $H_0: \delta \geq D_0$ versus $H_A: \delta < D_0$ is

$$\Phi\left(\frac{D_0 - \delta_1}{\frac{\sigma_d}{\sqrt{2n}}} + Z_\alpha\right)$$

where $\Phi()$ is the standard normal distribution function, $\delta_1$ is the actual value of the difference under the alternative hypothesis, and $Z_\alpha$ is the lower $\alpha$ percentile of the standard normal distribution. The sample size is determined using a binary search of possible values for $n$.

# Example 1 – Power Analysis

Suppose you want to consider the power of a balanced cross-over design with a binary endpoint where the test is computed based on the difference for sequence sample sizes between 50 and 200. The non-inferiority difference is -0.2, the actual difference is 0, and the estimated standard deviation of the paired differences is 1. The significance level is 0.05.

## 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 ......................................................**Power**
Higher Proportions Are .................................**Better**
Alpha............................................................**0.05**
n (Sample Size per Sequence) ......................**50 to 200 by 50**
D0 (Non-Inferiority Difference) .......................**-0.2**
D1 (Actual Difference)....................................**0**
Standard Deviation (SD) ................................**1**

## Output

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

**Numeric Results**
─────────────────────────────────────────────────────────────────────────────────────
Solve For:                    Power
Higher Proportions Are:       Better
Hypotheses:                   H0: Pt - Pc ≤ D0   vs.   H1: Pt - Pc > D0
─────────────────────────────────────────────────────────────────────────────────────

| | Sample Size | | Difference | | Standard | |
| | Sequence | Total | Non-Inferiority | Actual | Deviation | |
| Power | n | N | D0 | D1 | SD | Alpha |
|---|---|---|---|---|---|---|
| 0.63876 | 50 | 100 | -0.2 | 0 | 1 | 0.05 |
| 0.88171 | 100 | 200 | -0.2 | 0 | 1 | 0.05 |
| 0.96556 | 150 | 300 | -0.2 | 0 | 1 | 0.05 |
| 0.99074 | 200 | 400 | -0.2 | 0 | 1 | 0.05 |
─────────────────────────────────────────────────────────────────────────────────────

Power    The probability of rejecting a false null hypothesis when the alternative hypothesis is true.
n        The sample size in each sequence (or group).
N        The total sample size from both sequences. The sample is divided equally among sequences.
D0       The non-inferiority difference used to specify the hypothesis test.
D1       The actual difference at which power is calculated.
SD       The standard deviation of paired differences. This is estimated from a previous study.
Alpha    The probability of rejecting a true null hypothesis.

Non-Inferiority Tests for the Difference of Two Proportions in a 2x2 Cross-Over Design

**Summary Statements**
────────────────────────────────────────────────────────────

A 2×2 cross-over design will be used to test whether the treatment proportion is non-inferior to the standard proportion, with a non-inferiority margin of -0.2 (H0: Pt - Pc ≤ -0.2 versus H1: Pt - Pc > -0.2). The comparison will be made using a one-sided Z-test, with a Type I error rate (α) of 0.05. The standard deviation of paired differences is assumed to be 1. To detect a proportion difference (Pt - Pc) of 0 with a sample size of 50 in each sequence (totaling 100 subjects), the power is 0.63876.
────────────────────────────────────────────────────────────

**Dropout-Inflated Sample Size**
────────────────────────────────────────────────────────────

| Dropout Rate | Sample Size | | Dropout-Inflated Enrollment Sample Size | | Expected Number of Dropouts | |
|---|---|---|---|---|---|---|
| | n | N | n' | N' | d | D |
| 20% | 50 | 100 | 63 | 126 | 13 | 26 |
| 20% | 100 | 200 | 125 | 250 | 25 | 50 |
| 20% | 150 | 300 | 188 | 376 | 38 | 76 |
| 20% | 200 | 400 | 250 | 500 | 50 | 100 |

────────────────────────────────────────────────────────────

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.

n and N   The evaluable group and total sample sizes, respectively, at which power is computed (as entered by the user). If n subjects from each group are evaluated out of the n' subjects that are enrolled in the study, the design will achieve the stated power. N = 2n.

n' and N'   The number of subjects that should be enrolled in the study in order to obtain n and N evaluable subjects, based on the assumed dropout rate. n' is calculated by inflating n using the formula n' = n / (1 - DR), with n' always rounded up. (See Julious, S.A. (2010) pages 52-53, or Chow, S.C., Shao, J., Wang, H., and Lokhnygina, Y. (2018) pages 32-33.). N' = 2n'.

d and D   The expected number of group and total dropouts, respectively. d = n' - n and D = 2d.

**Dropout Summary Statements**
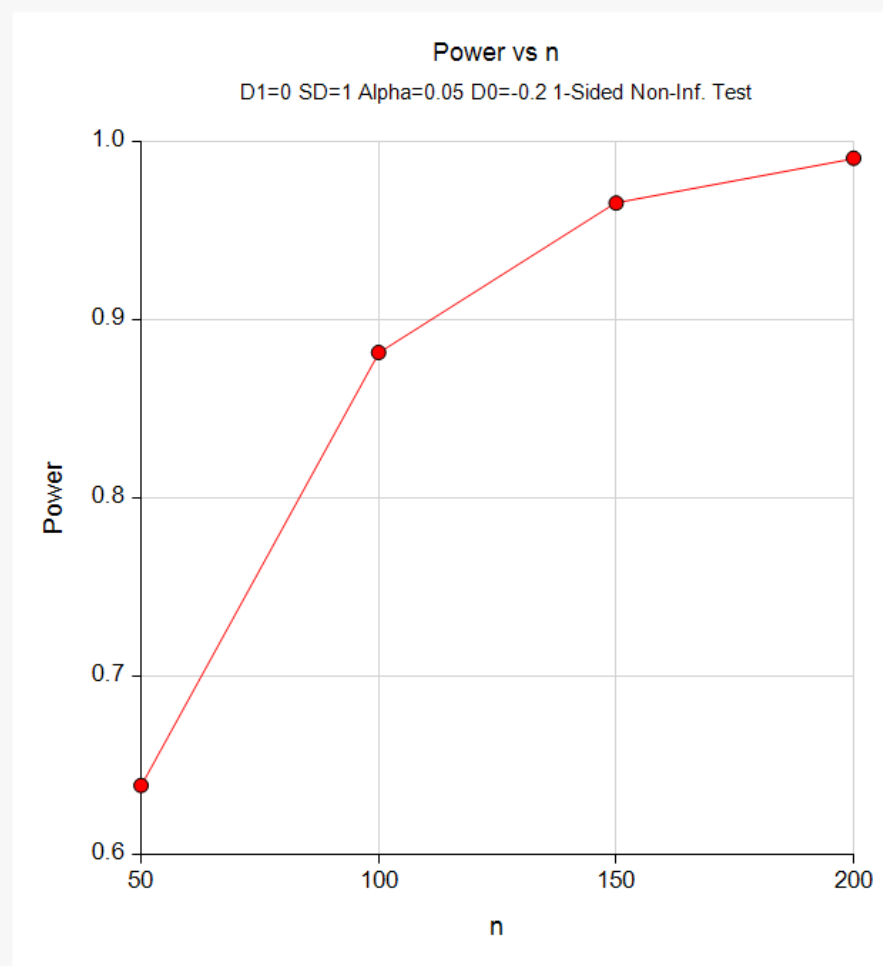────────────────────────────────────────────────────────────

Anticipating a 20% dropout rate, 63 subjects should be enrolled in each group to obtain final sample sizes of 50 subjects per group.
────────────────────────────────────────────────────────────

**References**
────────────────────────────────────────────────────────────

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

**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 100 per sequence is required for 90% power.

# Example 2 – Calculating Sample Size when Estimating the Standard Deviation from a Previous Study

This example demonstrates how to calculate the sample size when estimating the standard deviation of the paired differences from data in a previous study using the method in Chow, Shao, Wang, & Lokhnygina (2018) on pages 82 and 83. In this example we'll find the sample size required to detect a difference of 0 with 90% power in a test against a non-inferiority difference bound of -0.1 at a significance level of 0.05, with the standard deviation estimated using cell counts from a previous 2x2 cross-over study with equal sample size per sequence.

Assume that the following results were previously obtained from 280 subjects in a simple 2x2 cross-over trial comparing two inhalation devices, A and B. These results are similar to Table 3.2 of Lui (2016) on page 36 with a slight adjustment to sequence 1 that makes the sample sizes equal per sequence.

**SEQUENCE 1 (Control (A) → Treatment (B))**

|  |  | Period 2 (B) | | |
|---|---|---|---|---|
|  |  | Yes | No | Total |
| Period 1 (A) | Yes | 27 | 41 | 67 |
|  | No | 15 | 57 | 72 |
|  | Total | 42 | 98 | 140 |

**SEQUENCE 2 (Treatment (B) → Control (A))**

|  |  | Period 2 (A) | | |
|---|---|---|---|---|
|  |  | Yes | No | Total |
| Period 1 (B) | Yes | 38 | 16 | 54 |
|  | No | 32 | 54 | 86 |
|  | Total | 70 | 70 | 140 |

The paired differences of treatment – control for each subject within each sequence are

$$d_{ij} = x_{ijT} - x_{ijC}$$

$$= x_{ijB} - x_{ijA}$$

If $x_{ijk}$ is a binary variable, then $d_{ij}$ can take on the values 0, 1, and -1. If we summarize the results in sequence 1, then there are 15 subjects with $d_{ij} = 1$ (B = "Yes", A = "No"), 41 subjects with $d_{ij} = -1$ (B = "No", A = "Yes"), and 27 + 57 = 84 subjects with $d_{ij} = 0$ (B = "Yes", A = "Yes" or B = "No", A = "No"). The average paired difference, then, for sequence 1 is

$$\bar{d}_{1.} = \frac{1}{n}\sum_{j=1}^{n} d_{1j}$$

$$= \frac{15(1) + 41(-1) + 84(0)}{140}$$

$$= \frac{15 - 41}{140}$$

$$= -0.1857.$$

Similarly, the average paired difference from sequence 2 is

$$\bar{d}_{2.} = \frac{1}{n}\sum_{j=1}^{n} d_{2j}$$

$$= \frac{16(1) + 32(-1) + 92(0)}{140}$$

$$= \frac{16 - 32}{140}$$

$$= -0.1143.$$

The estimated overall treatment – control difference is

$$\hat{\delta} = \frac{1}{2n}\sum_{i=1}^{2}\sum_{j=1}^{n} d_{ij}$$

$$= \frac{\bar{d}_{1.} + \bar{d}_{2.}}{2}$$

$$= \frac{-0.1857 - 0.1143}{2}$$

$$= -0.15.$$

The estimated variance of paired differences, then, is

$$\hat{\sigma}_d^2 = \frac{1}{2(n-1)}\sum_{i=1}^{2}\sum_{j=1}^{n}\left(d_{ij} - \bar{d}_{i.}\right)^2$$

$$= \frac{15(1+0.1857)^2 + 41(-1+0.1857)^2 + 84(0+0.1857)^2 + 16(1+0.1143)^2 + 32(-1+0.1143)^2 + 92(0+0.1143)^2}{2(139)}$$

$$= 0.3502$$

with estimated standard deviation

$$\widehat{SD} = \hat{\sigma}_d$$

$$= \sqrt{\hat{\sigma}_d^2}$$

$$= \sqrt{0.3502}$$

$$= 0.5917.$$

# 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 Proportions Are ..................................**Better** |
| Power...............................................................**0.90** |
| Alpha................................................................**0.05** |
| D0 (Non-Inferiority Difference) .......................**-0.1** |
| D1 (Actual Difference)....................................**0** |
| Standard Deviation (SD) ................................**0.5917** |

# Output

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

**Numeric Results**

Solve For:                     Sample Size
Higher Proportions Are:        Better
Hypotheses:                    H0: Pt - Pc ≤ D0   vs.   H1: Pt - Pc > D0

| | Sample Size | | Difference | | Standard Deviation | |
|---|---|---|---|---|---|---|
| Power | Sequence n | Total N | Non-Inferiority D0 | Actual D1 | SD | Alpha |
| 0.90015 | 150 | 300 | -0.1 | 0 | 0.592 | 0.05 |

This report indicates that the required sample size for 90% power is 150 per sequence for a total of 300.

# Example 3 – Calculating Sample Size (Validation using Chow, Shao, Wang, & Lokhnygina (2018))

On page 86, Chow, Shao, Wang, & Lokhnygina (2018) presents an example of finding the sample size required to detect a difference of 0 in a test against a lower non-inferiority bound of -0.2 with 80% power and a significance level of 0.05 when the standard deviation of paired differences is 0.5. They compute the required sample size to be 20 per sequence.

## 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 3** 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 Proportions Are ..................................**Better**
Power...........................................................**0.8**
Alpha............................................................**0.05**
D0 (Non-Inferiority Difference) .......................**-0.2**
D1 (Actual Difference)....................................**0**
Standard Deviation (SD)................................**0.5**

## Output

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

**Numeric Results**
_____

Solve For:                    Sample Size
Higher Proportions Are:       Better
Hypotheses:                   H0: Pt - Pc ≤ D0   vs.   H1: Pt - Pc > D0
_____

|  | Sample Size | | Difference | | Standard | |
| Power | Sequence n | Total N | Non-Inferiority D0 | Actual D1 | Deviation SD | Alpha |
|---|---|---|---|---|---|---|
| 0.81191 | 20 | 40 | -0.2 | 0 | 0.5 | 0.05 |

_____

The result from **PASS** matches the result in Chow, Shao, Wang, & Lokhnygina (2018) exactly.