

Chapter 176

Superiority by a Margin 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 at 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, & 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 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 ($j = 1, \dots, n_i$) 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 $x_{ijk} = \bar{x}_{ijk} = \frac{1}{m} (\sum_{l=1}^m x_{ijkl})$, where x_{ijkl} represents the l^{th} binary response replicate ($l = 1, \dots, m$) from the j^{th} subject ($j = 1, \dots, n_i$) 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

$$\begin{aligned} d_{ij} &= x_{ijT} - x_{ijC} \\ &= x_{ij1} - x_{ij2} \end{aligned}$$

and the overall treatment – control difference as

$$\begin{aligned} \delta &= P_T - P_C \\ &= P_1 - P_2. \end{aligned}$$

Superiority by a Margin Tests for the Difference of Two Proportions in a 2x2 Cross-Over Design

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 σ_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

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

Superiority by a Margin Test Statistics

Higher Proportions Better

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

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

or equivalently

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

where D_0 is the superiority bound (i.e. the smallest difference ($P_T - P_C$) for which the treatment will be considered superior to the standard or control). When higher proportions are better, 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 α 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.

Superiority by a Margin Tests for the Difference of Two Proportions in a 2x2 Cross-Over Design

Higher Proportions Worse

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

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

or equivalently

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

where D_0 is the superiority bound (i.e. the largest difference ($P_T - P_C$) for which the treatment will be considered superior to the standard or control). When higher proportions are worse, 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 α if

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

where Z_α is the lower α percentile of the standard normal distribution.

Superiority by a Margin Power Calculations

Higher Proportions Better

According to Chow, Shao, Wang, & Lokhnygina (2018) page 84, the power for the one-sided superiority 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, δ_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 superiority 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, δ_1 is the actual value of the difference under the alternative hypothesis, and Z_α is the lower α percentile of the standard normal distribution. The sample size is determined using a binary search of possible values for n .

Procedure Options

This section describes the options that are specific to this procedure. These are located on the Design tab. For more information about the options of other tabs, go to the Procedure Window chapter.

Design Tab

The Design tab contains most of the parameters and options that you will be concerned with.

Solve For

Solve For

This option specifies the parameter to be calculated from the values of the other parameters. Under most conditions, you would select either *Power* or *Sample Size*.

Select *Sample Size* when you want to determine the sample size needed to achieve a given power and alpha level.

Select *Power* when you want to calculate the power of an experiment that has already been run.

Select *Effect Size (DI)* when you want to calculate the minimum effect size that can be detected for a particular design.

Test

Higher Proportions Are

Use this option to specify the direction of the test.

If Higher Proportions are “Better”, the alternative hypothesis is $H1: Pt - Pc > D0$.

If Higher Proportions are “Worse”, the alternative hypothesis is $H1: Pt - Pc < D0$.

Power and Alpha

Power

This option specifies one or more values for power. Power is the probability of rejecting a false null hypothesis, and is equal to one minus Beta. Beta is the probability of a type-II error, which occurs when a false null hypothesis is not rejected. In this procedure, a type-II error occurs when you fail to reject the null hypothesis of equal means when in fact the means are different.

Values must be between zero and one. Historically, the value of 0.80 (Beta = 0.20) was used for power. Now, 0.90 (Beta = 0.10) is also commonly used.

A single value may be entered here or a range of values such as *0.8 to 0.95 by 0.05* may be entered.

Alpha

This option specifies one or more values for the probability of a type-I error. A type-I error occurs when a true null hypothesis is rejected. In this procedure, a type-I error occurs when you reject the null hypothesis of equal means when in fact the means are equal.

Values must be between zero and one. Historically, the value of 0.05 has been used for alpha. This means that about one test in twenty will falsely reject the null hypothesis. You should pick a value for alpha that represents the risk of a type-I error you are willing to take in your experimental situation.

You may enter a range of values such as *0.01 0.05 0.10* or *0.01 to 0.10 by 0.01*.

Sample Size

n (Sample Size per Sequence)

This is the sample size of each sequence or group (AB and BA) in the cross-over design. The individual sequence sample sizes are assumed to be equal such that the total sample size is equal to $N = 2n$.

You can enter a single value such as 50 or a list of values using the syntax *50 100 150 200 250* or *50 to 250 by 50*.

Effect Size – Difference

D0 (Superiority Difference)

Specify the superiority difference.

When higher proportions are “Better”, the superiority difference is the smallest difference ($P_t - P_c$) for which the treatment will be considered superior to the standard or control.

When higher proportions are “Worse”, the superiority difference is the largest difference ($P_t - P_c$) for which the treatment will be considered superior to the standard or control.

You can enter a single value such as 0.1 or a series of values such as *0.1 0.15 0.2* or *0.1 to 0.2 by 0.05* in the range $-1 < D_0 < 1$, $D_0 \neq D_1$. When higher proportions are “Better”, D_0 should be greater than 0. When higher proportions are “Worse”, D_0 should be less than 0.

D1 (Actual Difference)

Enter a value for the actual difference ($P_t - P_c$) at which power is calculated. You can enter a single value such as 0.3 or a series of values such as *0.3 0.35 0.4* or *0.3 to 0.4 by 0.05* in the range $-1 < D_1 < 1$, $D_1 \neq D_0$. When higher proportions are “Better”, D_1 should be greater than D_0 . When higher proportions are “Worse”, D_1 should be less than D_0 .

Effect Size – Standard Deviation of Paired Differences

Standard Deviation (SD)

Enter a value for the standard deviation of the paired differences, SD.

Estimating SD using Previous 2x2 Cross-Over Data

The standard deviation may be estimated using cell counts from a previous 2x2 cross-over study with n subjects per sequence as described on pages 82 and 83 of Chow, Shao, Wang, & Lokhnygina (2018).

Assume x_{ij1} is the binary treatment response of the j th subject ($j = 1$ to n) in the i th sequence ($i = 1, 2$), and x_{ij2} is the binary control or reference response of the j th subject ($j = 1$ to n) in the i th sequence ($i = 1, 2$). Note that the number of subjects in sequence 1 is equal to the number of subjects in sequence 2 such that $n_1 = n_2 = n$.

Define

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

$$\bar{d}_{i} = (1/n)\sum_j d_{ij}$$

The formula for SD is then

$$SD = \sqrt{[(\sum_i \sum_j (d_{ij} - \bar{d}_{i})^2)/(2(n-1))]}.$$

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 superiority difference is 0.2, the actual difference is 0.4, and the estimated standard deviation of the paired differences is 1. The significance level is 0.05.

Setup

This section presents the values of each of the parameters needed to run this example. First, from the PASS Home window, load the **Superiority by a Margin Tests for the Difference of Two Proportions in a 2x2 Cross-Over Design** procedure window by expanding **Proportions**, then **Cross-Over (2x2) Design**, then clicking on **Superiority by a Margin**, and then clicking on **Superiority by a Margin Tests for the Difference of Two Proportions in a 2x2 Cross-Over Design**. You may then make the appropriate entries as listed below, or open **Example 1** by going to the **File** menu and choosing **Open Example Template**.

<u>Option</u>	<u>Value</u>
Design Tab	
Solve For	Power
Higher Proportions Are	Better
Alpha.....	0.05
n (Sample Size per Sequence).....	50 to 200 by 50
D0 (Superiority Difference)	0.2
D1 (Actual Difference)	0.4
Standard Deviation (SD).....	1

Annotated Output

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

Numeric Results

Numeric Results for a Superiority by a Margin Test

H0: $P_t - P_c \leq D_0$ vs. H1: $P_t - P_c > D_0$

	Sequence Sample Size n	Total Sample Size N	Superiority Difference D0	Actual Difference D1	Standard Deviation SD	Alpha
Power	50	100	0.200	0.400	1.000	0.050
	100	200	0.200	0.400	1.000	0.050
	150	300	0.200	0.400	1.000	0.050
	200	400	0.200	0.400	1.000	0.050

References

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

Report Definitions

Power is the probability of rejecting a false null hypothesis. It should be close to one.

n is the sample size in each sequence (or group).

N is the total sample size from both sequences. The sample is divided equally among sequences.

D0 is the superiority difference used to specify the hypothesis test.

D1 is the actual difference at which power is calculated.

SD is the standard deviation of paired differences. This is estimated from a previous study.

Alpha is the probability of rejecting a true null hypothesis. It should be small.

Superiority by a Margin Tests for the Difference of Two Proportions in a 2x2 Cross-Over Design

Summary Statements

For a 2x2 cross-over design, a sample size of 50 in each sequence for a total of 100 achieves 63.876% power to detect a difference of 0.400 using a one-sided superiority by a margin test against a bound of 0.200 with a significance level of 0.050 when the standard deviation of paired differences is 1.000.

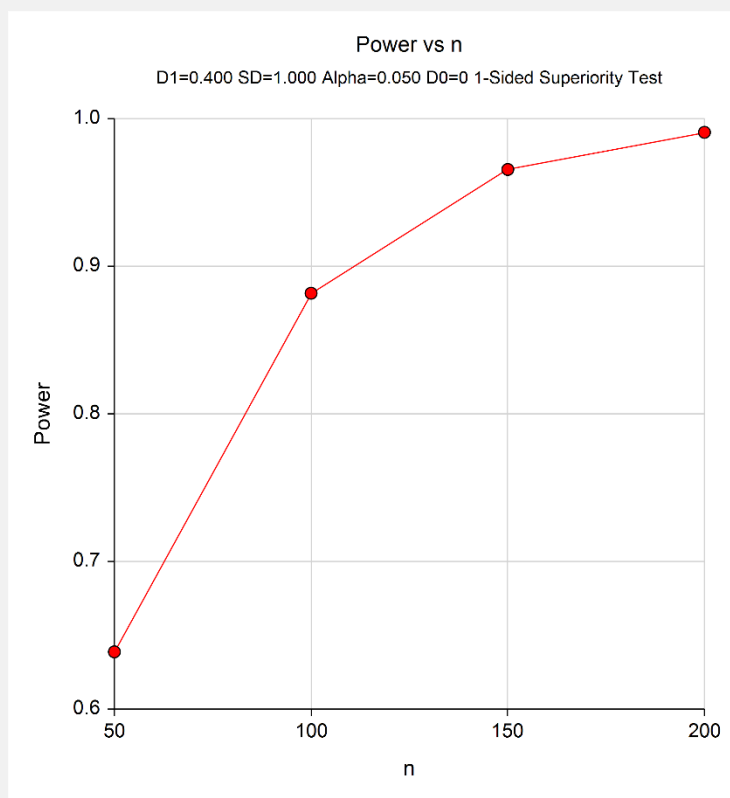
Dropout-Inflated Sample Size

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

Definitions

Dropout Rate (DR) is 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"). n and N are 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' are the number of subjects that should be enrolled in the study in order to end up with 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., and Wang, H. (2008) pages 39-40.). $N' = 2n'$. d and D are the expected number of group and total dropouts, respectively. $d = n' - n$ and $D = 2d$.

Charts Section



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 0.1 with 90% power in a test against a superiority 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))					SEQUENCE 2 (Treatment (B) → Control (A))				
		Period 2 (B)					Period 2 (A)		
		Yes	No	Total			Yes	No	Total
Period 1 (A)	Yes	27	41	67	Period 1 (B)	Yes	38	16	54
	No	15	57	72		No	32	54	86
Total		42	98	140	Total		70	70	140

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

$$\begin{aligned} d_{ij} &= x_{ijT} - x_{ijC} \\ &= x_{ijB} - x_{ijA} \end{aligned}$$

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

$$\begin{aligned} \bar{d}_{1.} &= \frac{1}{n} \sum_{j=1}^n d_{1j} \\ &= \frac{15(1) + 41(-1) + 84(0)}{140} \\ &= \frac{15 - 41}{140} \\ &= -0.1857. \end{aligned}$$

Superiority by a Margin Tests for the Difference of Two Proportions in a 2x2 Cross-Over Design

Similarly, the average paired difference from sequence 2 is

$$\begin{aligned}\bar{d}_2 &= \frac{1}{n} \sum_{j=1}^n d_{2j} \\ &= \frac{16(1) + 32(-1) + 92(0)}{140} \\ &= \frac{16 - 32}{140} \\ &= -0.1143.\end{aligned}$$

The estimated overall treatment – control difference is

$$\begin{aligned}\hat{\delta} &= \frac{1}{2n} \sum_{i=1}^2 \sum_{j=1}^n d_{ij} \\ &= \frac{\bar{d}_1 + \bar{d}_2}{2} \\ &= \frac{-0.1857 - 0.1143}{2} \\ &= -0.15.\end{aligned}$$

The estimated variance of paired differences, then, is

$$\begin{aligned}\hat{\sigma}_d^2 &= \frac{1}{2(n-1)} \sum_{i=1}^2 \sum_{j=1}^n (d_{ij} - \bar{d}_i)^2 \\ &= \frac{15(1 + 0.1857)^2 + 41(-1 + 0.1857)^2 + 84(0 + 0.1857)^2 + 16(1 + 0.1143)^2 + 32(-1 + 0.1143)^2 + 92(0 + 0.1143)^2}{2(139)} \\ &= 0.3502\end{aligned}$$

with estimated standard deviation

$$\begin{aligned}\widehat{SD} &= \hat{\sigma}_d \\ &= \sqrt{\hat{\sigma}_d^2} \\ &= \sqrt{0.3502} \\ &= 0.5917.\end{aligned}$$

Setup

This section presents the values of each of the parameters needed to run this example. First, from the PASS Home window, load the **Superiority by a Margin Tests for the Difference of Two Proportions in a 2x2 Cross-Over Design** procedure window by expanding **Proportions**, then **Cross-Over (2x2) Design**, then clicking on **Superiority by a Margin**, and then clicking on **Superiority by a Margin Tests for the Difference of Two Proportions in a 2x2 Cross-Over Design**. You may then make the appropriate entries as listed below, or open **Example 2** by going to the **File** menu and choosing **Open Example Template**.

<u>Option</u>	<u>Value</u>
Design Tab	
Solve For	Sample Size
Higher Proportions Are	Better
Power	0.90
Alpha	0.05
D0 (Superiority Difference)	0.1
D1 (Actual Difference)	0.2
Standard Deviation (SD)	0.5917

Output

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

Numeric Results

Numeric Results for a Superiority by a Margin Test						
H0: $P_t - P_c \leq D_0$ vs. H1: $P_t - P_c > D_0$						
	Sequence Sample Size	Total Sample Size	Superiority Difference	Actual Difference	Standard Deviation	Alpha
Power	n	N	D0	D1	SD	
0.90015	150	300	0.100	0.200	0.592	0.050

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.

We'll validate this procedure by computing the sample size for the same magnitude of difference between the actual difference and superiority difference--- 0.2. We'll set the superiority bound to 0.1 and the actual difference to 0.3 and should achieve the same sample size results.

Setup

This section presents the values of each of the parameters needed to run this example. First, from the PASS Home window, load the **Superiority by a Margin Tests for the Difference of Two Proportions in a 2x2 Cross-Over Design** procedure window by expanding **Proportions**, then **Cross-Over (2x2) Design**, then clicking on **Superiority by a Margin**, and then clicking on **Superiority by a Margin Tests for the Difference of Two Proportions in a 2x2 Cross-Over Design**. You may then make the appropriate entries as listed below, or open **Example 3** by going to the **File** menu and choosing **Open Example Template**.

<u>Option</u>	<u>Value</u>
Design Tab	
Solve For	Sample Size
Higher Proportions Are	Better
Power.....	0.80
Alpha.....	0.05
D0 (Superiority Difference)	0.1
D1 (Actual Difference)	0.3
Standard Deviation (SD).....	0.5

Output

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

Numeric Results

Numeric Results for a Superiority by a Margin Test						
H0: Pt - Pc ≤ D0 vs. H1: Pt - Pc > D0						
	Sequence Sample Size	Total Sample Size	Superiority Difference	Actual Difference	Standard Deviation	Alpha
Power	n	N	D0	D1	SD	Alpha
0.81191	20	40	0.100	0.300	0.500	0.050

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