

Chapter 328

Multi-Arm Non-Inferiority Tests for the Difference Between Treatment and Control Means Assuming Equal Variance

Introduction

This module computes power and sample size for multiple non-inferiority tests of treatment means versus a control mean. This chapter is based on the results in Machin, Campbell, Tan, and Tan (2018). In this design, there are k treatment groups and one control group. A mean is measured in each group. A total of k hypothesis tests are anticipated, each comparing a treatment group with the common control group using a non-inferiority t-test of the difference between two means.

The Bonferroni adjustment of the type I error rate may be optionally made because several comparisons are being tested using the same data. Making a multiplicity adjustment is usually recommended, but not always. In fact, Saville (1990) advocates not applying it and Machin, Campbell, Tan, and Tan (2018) include omitting it as a possibility.

Background

Whether you want to test several doses of a single treatment or several types of treatments, good research practice requires that each treatment be compared with a control. For example, a popular three-arm design consists of three groups: control, treatment A, and treatment B. Two non-inferiority tests are run: treatment A versus control and treatment B versus the same control. This design avoids having to obtain a second control group for treatment B. Besides the obvious efficiency in subjects, it may be easier to recruit subjects if their chances of receiving the new treatment are better than 50-50.

Technical Details

Suppose you want to compare k treatment groups with means μ_i and sample sizes N_i and one control group with mean μ_C and sample size N_C . The total sample size is $N = N_1 + N_2 + \cdots + N_k + N_C$.

Non-Inferiority Tests

A *non-inferiority test* tests that the treatment mean is not worse than the control mean by more than the non-inferiority margin (NIM). The actual direction of the hypothesis depends on the response variable being studied.

In the following sections, define $\delta_i = \mu_i - \mu_C$.

Case 1: High Values Good

In this case, higher response values are better. The hypotheses are arranged so that rejecting the null hypothesis implies that the treatment mean is no less than a small amount (NIM) below the control mean. The null and alternative hypotheses with are

$$H_{0i}: \mu_i - \mu_C \leq NIM \quad \text{vs.} \quad H_{1i}: \mu_i - \mu_C > NIM$$

$$H_{0i}: \mu_i \leq \mu_C + NIM \quad \text{vs.} \quad H_{1i}: \mu_i > \mu_C + NIM$$

$$H_{0i}: \delta_i \leq NIM \quad \text{vs.} \quad H_{1i}: \delta_i > NIM$$

where $NIM < 0$.

Case 2: High Values Bad

In this case, lower values are better. The hypotheses are arranged so that rejecting the null hypothesis implies that the treatment mean is no more than a small amount (NIM) above the control mean. The null and alternative hypotheses with are

$$H_{0i}: \mu_i - \mu_C \geq NIM \quad \text{vs.} \quad H_{1i}: \mu_i - \mu_C < NIM$$

$$H_{0i}: \mu_i \geq \mu_C + NIM \quad \text{vs.} \quad H_{1i}: \mu_i < \mu_C + NIM$$

$$H_{0i}: \delta_i \geq NIM \quad \text{vs.} \quad H_{1i}: \delta_i < NIM$$

where $NIM > 0$.

Two-Sample Equal-Variance T-Test Statistic

Under the null hypothesis, this test assumes that the two groups of data are simple random samples from a single population of normally distributed values that all have the same mean and variance. This assumption implies that the data are continuous, and their distribution is symmetric. The calculation of the test statistic for the case when higher response values are better is as follows.

$$t_{df} = \frac{(\bar{x}_i - \bar{x}_c) - NIM}{\sqrt{\frac{(N_i - 1)s_i^2 + (N_c - 1)s_c^2}{N_i + N_c - 2} \left(\frac{1}{N_i} + \frac{1}{N_c} \right)}}$$

where

$$\bar{X}_i = \frac{\sum_{j=1}^{N_i} X_{ij}}{N_i}$$

$$s_i = \sqrt{\left(\frac{\sum_{j=1}^{N_i} (X_{ij} - \bar{X}_i)^2}{(N_i - 1)} \right)}$$

$$df = N_i + N_c - 2$$

This t -statistic follows a t distribution with $N_i + N_c - 2$ degrees of freedom.

Power Calculation

The power of this test is computed using the noncentral t distribution with $N_i + N_c - 2$ degrees of freedom and non-centrality parameter

$$\lambda = \frac{\mu_i - \mu_c - NIM}{\sigma \sqrt{\frac{1}{N_i} + \frac{1}{N_c}}}$$

Multiplicity Adjustment

Because k t-tests between treatment groups and the control group are run when analyzing the results of this study, many statisticians recommend that the Bonferroni adjustment be applied. This adjustment is easy to apply: the value of alpha that is used in the test is found by dividing the original alpha by the number of tests. For example, if the original alpha is set at 0.05 and the number of treatment (not including the control) groups is five, the individual tests will be conducted using an alpha of 0.01.

The main criticism of this procedure is that if there are many tests, the value of alpha becomes very small. To mitigate against this complaint, some statisticians recommend separating the treatment groups into those that are of primary interest and those that are of secondary interest. The Bonferroni adjustment is made by using the number of primary treatments rather than the total number of treatments.

There are some who advocate ignoring the adjustment entirely in the case of randomized clinical trials. See for example Saville (1990) and the discussion in chapter 14 of Machin, Campbell, Tan, and Tan (2018).

Size of the Control Group

Because the control group is used over and over, some advocate increasing the number of subjects in this group. The standard adjustment is to include \sqrt{k} subjects in the control group for each subject in one of the treatment groups. See Machin, Campbell, Tan, and Tan (2018, pages 231-232). Note that often, the treatment groups all have the same size.

Example 1 – Finding the Sample Size

A parallel-group clinical trial is being designed to compare three treatment therapies against the standard therapy. Higher values of the response are desirable. Suppose the standard therapy has a mean response of 9.3 with a standard deviation of 2.5. The investigators would like a sample size large enough to find statistical significance at the 0.05 level if the actual mean responses of the three treatments are 9.1, 9.2 and 9.3, the power of each test is 0.80, and the non-inferiority margin is -10% of 9.3 = -0.93. They want to consider a range of standard deviations from 2.0 to 3.0.

Following standard procedure, the control group multiplier will be set to $\sqrt{k} = \sqrt{3} = 1.732$ since the control group is used for three comparisons in this design.

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	Sample Size
Higher Means Are	Better (H1: $\delta > \text{NIM}$)
Power of Each Test	0.8
Overall Alpha	0.05
Bonferroni Adjustment	Standard Bonferroni
Group Allocation	Enter Group Allocation Pattern, solve for group sample sizes
NIM (Non-Inferiority Margin)	-0.93
Control Mean	9.3
Control Sample Size Allocation	1.732
Set A Number of Groups	1
Set A Mean	9.1
Set A Sample Size Allocation	1
Set B Number of Groups	1
Set B Mean	9.2
Set B Sample Size Allocation	1
Set C Number of Groups	1
Set C Mean	9.3
Set C Sample Size Allocation	1
Set D Number of Groups	0
More	Unchecked
σ (Standard Deviation)	2 2.5 3

Multi-Arm Non-Inferiority Tests for the Difference Between Treatment and Control Means Assuming Equal Variance

Output

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

Numeric Reports

Numeric Results

Solve For: [Sample Size](#)
 Group Allocation: Enter Group Allocation Pattern, solve for group sample sizes
 Test Type: T-Test
 Higher Means Are: Better
 Hypotheses: $H_0: \delta \leq \text{NIM}$ vs. $H_1: \delta > \text{NIM}$
 Number of Groups: 4
 Bonferroni Adjustment: Standard Bonferroni (Divisor = 3)

Comparison	Power		Sample Size		Mean μ_i	Difference δ_i	Non-Inferiority Margin NIM	Standard Deviation σ	Alpha	
	Target	Actual	Ni	Allocation					Overall	Bonferroni-Adjusted
Control			184	1.732	9.3			2.0		
vs A	0.8	0.80331	106	1.000	9.1	-0.2	-0.93	2.0	0.05	0.01667
vs B	0.8	0.89651	106	1.000	9.2	-0.1	-0.93	2.0	0.05	0.01667
vs C	0.8	0.95258	106	1.000	9.3	0.0	-0.93	2.0	0.05	0.01667
Total			502							
Control			284	1.732	9.3			2.5		
vs A	0.8	0.80002	164	1.000	9.1	-0.2	-0.93	2.5	0.05	0.01667
vs B	0.8	0.89408	164	1.000	9.2	-0.1	-0.93	2.5	0.05	0.01667
vs C	0.8	0.95107	164	1.000	9.3	0.0	-0.93	2.5	0.05	0.01667
Total			776							
Control			409	1.732	9.3			3.0		
vs A	0.8	0.80179	236	1.000	9.1	-0.2	-0.93	3.0	0.05	0.01667
vs B	0.8	0.89527	236	1.000	9.2	-0.1	-0.93	3.0	0.05	0.01667
vs C	0.8	0.95176	236	1.000	9.3	0.0	-0.93	3.0	0.05	0.01667
Total			1117							

Comparison	The group that is involved in the comparison between the treatment and control displayed on this report line. The comparison is made using the difference.
Target Power	The power desired. Power is probability of rejecting a false null hypothesis for this comparison. This power is of the comparison shown on this line only.
Actual Power	The power actually achieved.
Ni	The number of subjects in the ith group. The total sample size shown below the groups is equal to the sum of all individual group sample sizes.
Allocation	The group sample size allocation ratio of the ith group. The value on each row represents the relative number of subjects assigned to the group.
μ_i	The mean of the ith group at which the power is computed. The first row contains μ_c , the control group mean.
δ_i	The difference between the ith treatment mean and the control mean ($\mu_i - \mu_c$) at which the power is computed.
σ	The standard deviation of the responses within each group.
NIM	The margin of non-inferiority in the scale of the mean difference. $\text{NIM} < 0$.
Overall Alpha	The probability of rejecting at least one of the comparisons in this experiment when each null hypothesis is true.
Bonferroni Alpha	The adjusted significance level at which each individual comparison is made.

Summary Statements

A parallel, 4-group design (with one control group and 3 treatment groups) will be used to test whether the mean for each treatment group is non-inferior to the control group mean, with a non-inferiority margin of -0.93 ($H_0: \delta \leq -0.93$ versus $H_1: \delta > -0.93$, $\delta = \mu_i - \mu_c$). In this study, higher means are considered to be better. The non-inferiority hypotheses will be evaluated using 3 one-sided, two-sample, Bonferroni-adjusted, equal-variance t-tests, with an overall (experiment-wise) Type I error rate (α) of 0.05. The common standard deviation for all groups is assumed to be 2. The control group mean is assumed to be 9.3. To detect the treatment means 9.1, 9.2, and 9.3 with at least 80% power for each test, the control group sample size needed will be 184 and the number of needed subjects for the treatment groups will be 106, 106, and 106 (totaling 502 subjects overall).

Multi-Arm Non-Inferiority Tests for the Difference Between Treatment and Control Means Assuming Equal Variance

Dropout-Inflated Sample Size

Group	Dropout Rate	Sample Size Ni	Dropout- Inflated Enrollment Sample Size Ni'	Expected Number of Dropouts Di
1	20%	184	230	46
2	20%	106	133	27
3	20%	106	133	27
4	20%	106	133	27
Total		502	629	127
1	20%	284	355	71
2	20%	164	205	41
3	20%	164	205	41
4	20%	164	205	41
Total		776	970	194
1	20%	409	512	103
2	20%	236	295	59
3	20%	236	295	59
4	20%	236	295	59
Total		1117	1397	280

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 Lokhnygina, 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 230, 133, 133, and 133 subjects should be enrolled to obtain final group sample sizes of 184, 106, 106, and 106 subjects.

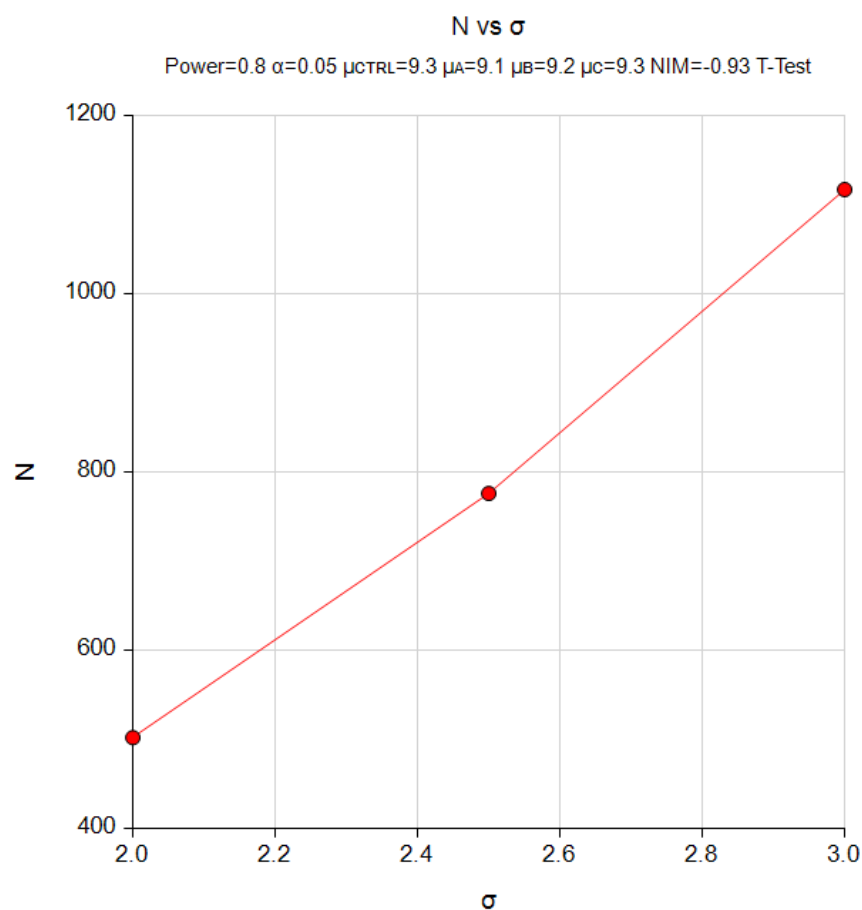
References

- Blackwelder, W.C. 1998. 'Equivalence Trials.' In Encyclopedia of Biostatistics, John Wiley and Sons. New York. Volume 2, 1367-1372.
- Chow, S.C., Shao, J., Wang, H., and Lokhnygina, Y. 2018. Sample Size Calculations in Clinical Research, 3rd Edition. Chapman & Hall/CRC. Boca Raton, FL. Pages 86-88.
- Julious, Steven A. 2004. 'Tutorial in Biostatistics. Sample sizes for clinical trials with Normal data.' Statistics in Medicine, 23:1921-1986.
- Machin, D., Campbell, M.J., Tan, S.B, and Tan, S.H. 2018. Sample Sizes for Clinical, Laboratory, and Epidemiology Studies, 4th Edition. Wiley Blackwell.

This report shows the numeric results of this power study. Notice that the results are shown in blocks of three rows at a time. Each block represents a single design.

Plots Section

Plots



This plot gives a visual presentation to the results in the Numeric Report. We can quickly see the impact on the sample size of changing the standard deviation.

Example 2 – Validation using a Previously Validated Procedure

We could not find a validation result in the statistical literature, so we will use a previously validated **PASS** procedure (**Two-Sample T-Tests for Non-Inferiority Assuming Equal Variance**) to produce the results for the following example.

A parallel-group clinical trial is being designed to compare three treatment therapies against the standard therapy. Higher values of the response are desirable. Suppose the standard therapy has a mean response of 9.3 with a standard deviation of 2.5. The investigators would like a sample size large enough to find statistical significance at the 0.05 level if the actual mean responses of the three treatments are 9.1, 9.2 and 9.3, the power of each test is 0.80, and the non-inferiority margin is -10% of 9.3 = -0.93.

The sample sizes of all groups will be equal.

The **Two-Sample T-Tests for Non-Inferiority Assuming Equal Variance** procedure is set up as follows.

Design Tab

Solve For **Sample Size**
 Higher Means Are **Better (H1: $\delta > -NIM$)**
 Power **0.8**
 Alpha **0.016667** (which is Alpha / k)
 Group Allocation **Equal (N1 = N2)**
 NIM (Non-Inferiority Margin) **-0.93**
 δ (Actual Difference to Detect) **-0.2 -0.1 0**
 σ (Standard Deviation) **2.5**

This set of options generates the following report.

Numeric Results

Solve For: [Sample Size](#)
 Test Type: Two-Sample Equal-Variance T-Test
 Difference: $\delta = \mu_1 - \mu_2 = \mu_T - \mu_R$
 Higher Means Are: Better
 Hypotheses: $H_0: \delta \leq -NIM$ vs. $H_1: \delta > -NIM$

Power		Sample Size			Non-Inferiority Margin -NIM	Mean Difference δ	Standard Deviation σ	Alpha
Target	Actual	N1	N2	N				
0.8	0.80001	208	208	416	-0.93	-0.2	2.5	0.01667
0.8	0.80218	162	162	324	-0.93	-0.1	2.5	0.01667
0.8	0.80132	129	129	258	-0.93	0.0	2.5	0.01667

In order to maintain a power of 80% for all three groups, it is apparent that the groups will all need to have a sample size of 208. This table contains the validation values. We will now run these values through the current procedure and compare the results with these values.

Multi-Arm Non-Inferiority Tests for the Difference Between Treatment and Control Means Assuming Equal Variance

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 (H1: $\delta > \text{NIM}$)**
 Power of Each Test **0.8**
 Overall Alpha **0.05**
 Bonferroni Adjustment **Standard Bonferroni**
 Group Allocation **Equal (Nc = N1 = N2 = ...)**
 NIM (Non-Inferiority Margin) **-0.93**
 Control Mean **9.3**
 Set A Number of Groups..... **1**
 Set A Mean **9.1**
 Set B Number of Groups..... **1**
 Set B Mean **9.2**
 Set C Number of Groups **1**
 Set C Mean..... **9.3**
 Set D Number of Groups **0**
 More..... **Unchecked**
 σ (Standard Deviation)..... **2.5**

Output

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

Numeric Results

Solve For: [Sample Size](#)
 Group Allocation: Equal (Nc = N1 = N2 = ...)
 Test Type: T-Test
 Higher Means Are: Better
 Hypotheses: H0: $\delta \leq \text{NIM}$ vs. H1: $\delta > \text{NIM}$
 Number of Groups: 4
 Bonferroni Adjustment: Standard Bonferroni (Divisor = 3)

Comparison	Power		Sample Size Ni	Mean μ_i	Difference δ_i	Non-Inferiority Margin NIM	Standard Deviation σ	Alpha	
	Target	Actual						Overall	Bonferroni- Adjusted
Control			208	9.3			2.5		
vs A	0.8	0.80000	208	9.1	-0.2	-0.93	2.5	0.05	0.01667
vs B	0.8	0.89406	208	9.2	-0.1	-0.93	2.5	0.05	0.01667
vs C	0.8	0.95106	208	9.3	0.0	-0.93	2.5	0.05	0.01667
Total			832						

As you can see, the sample sizes are all 208, which match the largest sample size found in the validation run above. The procedure is validated.