

## Chapter 339

# Multi-Arm Tests for the Difference Between Treatment and Control Means Allowing Unequal Variance

---

## Introduction

This module computes power and sample size for multiple comparisons of treatment means versus a control mean when no assumption of equal variances for the group populations is made. This is commonly known as the Aspin-Welch test, Welch's t-test (Welch, 1937), or the Satterthwaite method. 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 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 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  $\mu_i$  and sample sizes  $N_i$  and one control group with mean  $\mu_C$  and sample size  $N_C$ . The total sample size is  $N = N_1 + N_2 + \cdots + N_k + N_C$ .

The statistical hypotheses for two-sided tests are

$$H_{0i}: \mu_i = \mu_C \text{ vs. } H_{1i}: \mu_i \neq \mu_C$$

and for one-sided tests the hypotheses are

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

or

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

If we define  $\delta_i = \mu_i - \mu_C$ , these are equivalent to

$$H_{0i}: \delta_i = 0 \text{ vs. } H_{1i}: \delta_i \neq 0 \text{ for } i = 1, 2, \dots, k$$

$$H_{0i}: \delta_i \leq 0 \text{ vs. } H_{1i}: \delta_i > 0 \text{ for } i = 1, 2, \dots, k$$

$$H_{0i}: \delta_i \geq 0 \text{ vs. } H_{1i}: \delta_i < 0 \text{ for } i = 1, 2, \dots, k$$

For convenience, these hypotheses are collectively referred to as

$$H_0: \delta = 0 \text{ vs. } H_1: \delta \neq 0$$

$$H_0: \delta \leq 0 \text{ vs. } H_1: \delta > 0$$

$$H_0: \delta \geq 0 \text{ vs. } H_1: \delta < 0$$

## Test Statistic

A suitable Type I error probability is chosen for the test, the data are collected, and a  $t$ -statistic is generated using the formula

$$t = \frac{\bar{x}_i - \bar{x}_C}{\sqrt{\frac{s_i^2}{N_i} + \frac{s_C^2}{N_C}}}$$

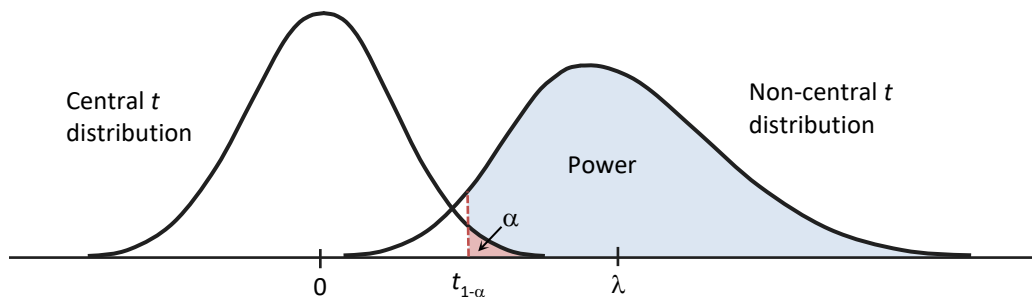
This  $t$ -statistic follows a  $t$  distribution approximately, with estimated degrees of freedom

$$df' = \frac{\left(\frac{s_i^2}{N_i} + \frac{s_C^2}{N_C}\right)^2}{\frac{1}{N_i - 1} \left(\frac{s_i^2}{N_i}\right)^2 + \frac{1}{N_C - 1} \left(\frac{s_C^2}{N_C}\right)^2}$$

## Power Calculation

This section describes the procedure for computing the power from the  $N_i$  and  $N_C$ ,  $\alpha$ , the assumed  $\mu_i$  and  $\mu_C$ , and the assumed standard deviations,  $\sigma_i$  and  $\sigma_C$ . Two good references for these general methods are Julious (2010) and Chow, Shao, Wang, and Lokhnygina (2018), although these texts do not specifically cover the Aspin-Welch-Satterthwaite t-test methods.

The figure below gives a visual representation for the calculation of power for a one-sided test.



If we call the assumed difference between the means,  $\delta_i = \mu_i - \mu_C$ , the steps for calculating the power are as follows:

1. Find  $t_{1-\alpha}$  based on the central- $t$  distribution with degrees of freedom,

$$df_i = \frac{\left(\frac{\sigma_i^2}{N_i} + \frac{\sigma_C^2}{N_C}\right)^2}{\frac{1}{N_i - 1} \left(\frac{\sigma_i^2}{N_i}\right)^2 + \frac{1}{N_C - 1} \left(\frac{\sigma_C^2}{N_C}\right)^2}.$$

2. Calculate the non-centrality parameter:

$$\lambda_i = \frac{\delta_i}{\sqrt{\frac{\sigma_i^2}{N_i} + \frac{\sigma_C^2}{N_C}}}.$$

3. Calculate the power as the probability that the test statistic  $t$  is greater than  $t_{1-\alpha}$  under the non-central- $t$  distribution with non-centrality parameter  $\lambda_i$ :

$$Power = \Pr_{Non-central-t}(t > t_{1-\alpha} | df_i, \lambda_i).$$

The algorithms for calculating power for the opposite direction and the two-sided hypotheses are analogous to this method.

When solving for something other than power, **PASS** uses this same power calculation formulation, but performs a search to determine that parameter.

## 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. Suppose the standard therapy has a mean response of 9.3. The investigators would like a sample size large enough to find statistical significance at the 0.05 level when the actual mean responses of the three treatments are all 7.6 and the power is 0.80 in each test. They want to consider standard deviations of 2.7 in the control group and 2.1 in the three treatment groups. They want to see what happens if the standard deviations are multiplied by 0.8, 1.0, and 1.2. The tests will be two-sided.

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 .....	<b>Sample Size</b>
Alternative Hypothesis .....	<b>Two-Sided (H1: <math>\delta \neq 0</math>)</b>
Power of Each Test .....	<b>0.80</b>
Overall Alpha .....	<b>0.05</b>
Bonferroni Adjustment .....	<b>Standard Bonferroni</b>
Group Allocation .....	<b>Enter Group Allocation Pattern, solve for group sample sizes</b>
Control Mean .....	<b>9.3</b>
Control Standard Deviation.....	<b>2.7</b>
Control Sample Size Allocation.....	<b>1.732</b>
Set A Number of Groups.....	<b>3</b>
Set A Mean .....	<b>7.6</b>
Set A Standard Deviation.....	<b>2.1</b>
Set A Sample Size Allocation .....	<b>1</b>
Set B Number of Groups.....	<b>0</b>
Set C Number of Groups .....	<b>0</b>
Set D Number of Groups .....	<b>0</b>
More.....	<b>Unchecked</b>
Add sets of standard deviations with .....	<b>Checked</b>
different magnitudes, but identical	
ratio patterns	
K ( $\sigma$ Multiplier) .....	<b>0.8 1 1.2</b>

## Multi-Arm Tests for the Difference Between Treatment and Control Means Allowing Unequal 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: Unequal-Variance T-Test  
 Hypotheses:  $H_0: \delta = 0$  vs.  $H_1: \delta \neq 0$   
 Number of Groups: 4  
 Bonferroni Adjustment: Standard Bonferroni (Divisor = 3)

Comparison	Power		Sample Size		Mean		Standard Deviation		Alpha	
	Target	Actual	Ni	Allocation	Value $\mu_i$	Difference $\delta_i$	Value $\sigma_i$	Multiplier K	Overall	Bonferroni-Adjusted
Control			38	1.732	9.3		2.16	0.8		
vs A1	0.8	0.81761	22	1.000	7.6	-1.7	1.68	0.8	0.05	0.01667
vs A2	0.8	0.81761	22	1.000	7.6	-1.7	1.68	0.8	0.05	0.01667
vs A3	0.8	0.81761	22	1.000	7.6	-1.7	1.68	0.8	0.05	0.01667
Total			104							
Control			57	1.732	9.3		2.70	1.0		
vs A1	0.8	0.80806	33	1.000	7.6	-1.7	2.10	1.0	0.05	0.01667
vs A2	0.8	0.80806	33	1.000	7.6	-1.7	2.10	1.0	0.05	0.01667
vs A3	0.8	0.80806	33	1.000	7.6	-1.7	2.10	1.0	0.05	0.01667
Total			156							
Control			81	1.732	9.3		3.24	1.2		
vs A1	0.8	0.80759	47	1.000	7.6	-1.7	2.52	1.2	0.05	0.01667
vs A2	0.8	0.80759	47	1.000	7.6	-1.7	2.52	1.2	0.05	0.01667
vs A3	0.8	0.80759	47	1.000	7.6	-1.7	2.52	1.2	0.05	0.01667
Total			222							

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.
$\mu_i$	The mean of the ith group at which the power is computed. The first row contains $\mu_c$ , the control group mean.
$\delta_i$	The difference between the ith treatment mean and the control mean ( $\mu_i - \mu_c$ ) at which the power is computed.
$\sigma_i$	The standard deviation of the responses within this group.
K	The multiplier that was applied to form the group standard deviations shown on this line.
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 different from the control group mean ( $H_0: \delta = 0$  versus  $H_1: \delta \neq 0$ ,  $\delta = \mu_i - \mu_c$ ). The hypotheses will be evaluated using 3 two-sided, two-sample, Bonferroni-adjusted, unequal-variance (Welch's) t-tests, with an overall (experiment-wise) Type I error rate ( $\alpha$ ) of 0.05. The group standard deviations (beginning with the control group) are assumed to be 2.16, 1.68, 1.68, and 1.68. The control group mean is assumed to be 9.3. To detect the treatment means 7.6, 7.6, and 7.6 with at least 80% power for each test, the control group sample size needed will be 38 and the number of needed subjects for the treatment groups will be 22, 22, and 22 (totaling 104 subjects overall).

## Multi-Arm Tests for the Difference Between Treatment and Control Means Allowing Unequal Variance

## Dropout-Inflated Sample Size

Group	Dropout Rate	Sample Size Ni	Dropout- Inflated Enrollment Sample Size Ni'	Expected Number of Dropouts Di
1	20%	38	48	10
2	20%	22	28	6
3	20%	22	28	6
4	20%	22	28	6
Total		104	132	28
1	20%	57	72	15
2	20%	33	42	9
3	20%	33	42	9
4	20%	33	42	9
Total		156	198	42
1	20%	81	102	21
2	20%	47	59	12
3	20%	47	59	12
4	20%	47	59	12
Total		222	279	57

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 48, 28, 28, and 28 subjects should be enrolled to obtain final group sample sizes of 38, 22, 22, and 22 subjects.

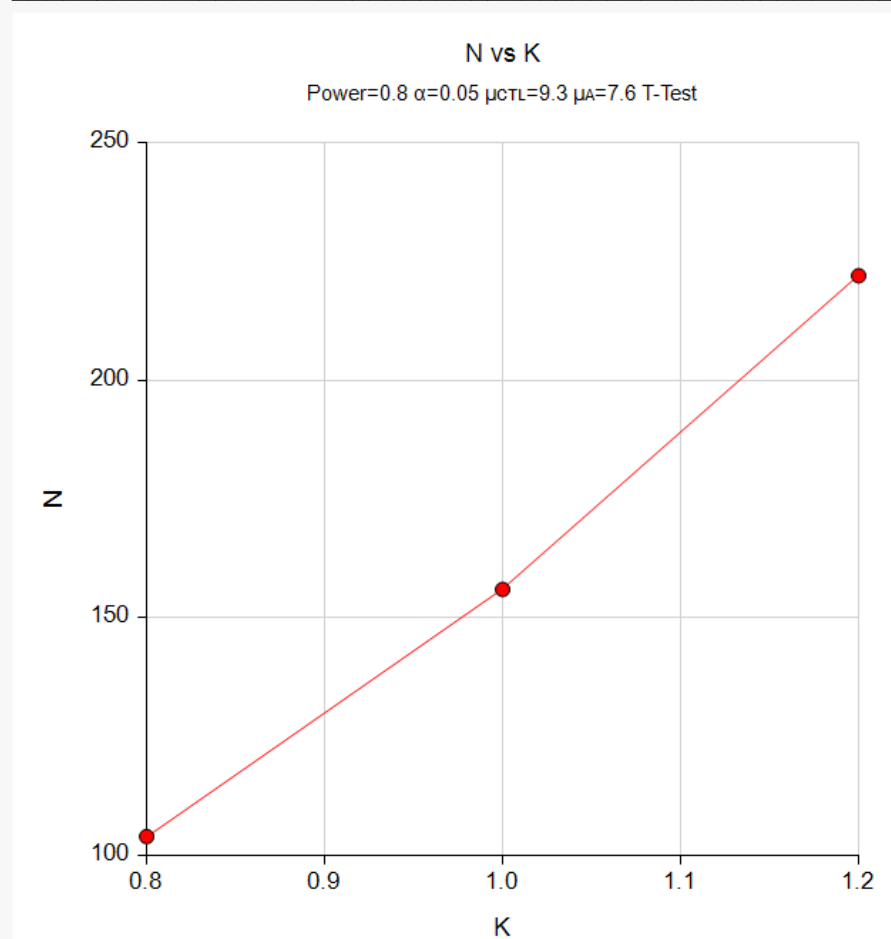
## References

- Chow, S.C., Shao, J., Wang, H., and Lohknygina, Y. 2018. Sample Size Calculations in Clinical Research, 3rd Edition. Chapman & Hall/CRC. Boca Raton, FL. Pages 86-88.
- 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.
- Julious, Steven A. 2004. 'Tutorial in Biostatistics. Sample sizes for clinical trials with Normal data.' Statistics in Medicine, 23:1921-1986.
- Welch, B.L. 1938. 'The significance of the difference between two means when the population variances are unequal.' Biometrika, 29, 350-362.
- Zar, Jerrold H. 1984. Biostatistical Analysis (Second Edition). Prentice-Hall. Englewood Cliffs, New Jersey.

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 magnitude.



## 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 Allowing Unequal 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. Suppose the standard therapy has a mean response of 9.3. The investigators would like a sample size large enough to find statistical significance at the 0.05 level when the actual mean responses of the three treatments are all 7.6 and the power is 0.80 in each test. They want to consider standard deviations of 2.7 in the control group and 2.1 in the three treatment groups. The tests will be two-sided.

The sample sizes of all groups will be equal.

The **Two-Sample T-Tests Allowing Unequal Variance** procedure is set up as follows.

### Design Tab

Solve For ..... **Sample Size**  
 Alternative Hypothesis ..... **Two-Sided**  
 Power..... **0.8**  
 Alpha..... **0.01667** (which is Alpha / k)  
 Group Allocation ..... **Equal (N1 = N2)**  
 Input Type..... **Means**  
 $\mu_1$ ..... **9.3**  
 $\mu_2$ ..... **7.6**  
 $\sigma_1$  ..... **2.7**  
 $\sigma_2$  ..... **2.1**

This set of options generates the following report.

### Numeric Results

Solve For: **Sample Size**  
 Test Type: Two-Sample Welch's Unequal-Variance T-Test  
 Difference:  $\delta = \mu_1 - \mu_2$   
 Hypotheses:  $H_0: \delta = 0$  vs.  $H_1: \delta \neq 0$

Power		Sample Size			Mean			Standard Deviation		Alpha
Target	Actual	N1	N2	N	$\mu_1$	$\mu_2$	Difference $\delta$	$\sigma_1$	$\sigma_2$	
0.8	0.80076	44	44	88	9.3	7.6	1.7	2.7	2.1	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 44 per group. 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 Tests for the Difference Between Treatment and Control Means Allowing Unequal 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**  
 Alternative Hypothesis ..... **Two-Sided ( $H_1: \delta \neq 0$ )**  
 Power of Each Test ..... **0.80**  
 Overall Alpha ..... **0.05**  
 Bonferroni Adjustment ..... **Standard Bonferroni**  
 Group Allocation ..... **Equal ( $N_c = N_1 = N_2 = \dots$ )**  
 Control Mean ..... **9.3**  
 Control Standard Deviation ..... **2.7**  
 Set A Number of Groups ..... **3**  
 Set A Mean ..... **7.6**  
 Set A Standard Deviation ..... **2.1**  
 Set B Number of Groups ..... **0**  
 Set C Number of Groups ..... **0**  
 Set D Number of Groups ..... **0**  
 More ..... **Unchecked**  
 Add sets of standard deviations with ..... **Unchecked**  
     different magnitudes, but identical  
     ratio patterns

## Output

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

### Numeric Results

Solve For: **Sample Size**  
 Group Allocation: **Equal ( $N_c = N_1 = N_2 = \dots$ )**  
 Test Type: **Unequal-Variance T-Test**  
 Hypotheses:  **$H_0: \delta = 0$  vs.  $H_1: \delta \neq 0$**   
 Number of Groups: **4**  
 Bonferroni Adjustment: **Standard Bonferroni (Divisor = 3)**

Comparison	Power		Sample Size $N_i$	Mean		Standard Deviation $\sigma_i$	Alpha	
	Target	Actual		Value $\mu_i$	Difference $\delta_i$		Overall	Bonferroni- Adjusted
Control			44	9.3		2.7		
vs A1	0.8	0.80073	44	7.6	-1.7	2.1	0.05	0.01667
vs A2	0.8	0.80073	44	7.6	-1.7	2.1	0.05	0.01667
vs A3	0.8	0.80073	44	7.6	-1.7	2.1	0.05	0.01667
Total			176					

As you can see, the sample sizes are all 44 which match the largest sample size found in the validation run above. The procedure is validated.