

# Module 3: Power

What is Power? + Calculating Power.

# Motivation

---

We haven't yet answered... How much data should we collect?

# Outcomes of hypothesis testing

	$H_0$ is True	$H_0$ is False
Reject $H_0$	Type I ( _____ )	No Error ( _____ )
Fail to Reject $H_0$	No Error ( _____ )	Type II ( _____ )

# What power means

## **Power** $(1 - \beta)$

The **power** of a test is the probability of correctly rejecting the null hypothesis if a particular alternative scenario is true.

$\text{Power} = 1 - \beta = P(\text{reject } H_0 \mid \text{a meaningful difference does in fact exist})$

# Why $\beta$ is harder to find than $\alpha$

---

Power depends on:

- Variability ( $\sigma^2$ )
- Sample size (replications)
- Design choices

**A key idea:** You cannot compute power for “something differs”

# Example 3.1: Shoe Type (extended)

Suppose, we are preparing to conduct a study for:

- $t = 4$  shoe types
- $r = 3$  runners per shoe
- CRD randomization

# Meaningful difference

$\delta$  = difference worth detecting

- Example: Fastest vs slowest shoe differs by **10 seconds**

# What we must specify

## **i Information Needed to Calculate the Power of a Test**

1. a specific alternative (e.g.,  $\delta$ )
2. replications per treatment ( $r$ )
3. variability estimate ( $\sigma^2$ )
4. significance level ( $\alpha$ )



# Reminder: the ANOVA test

---

$$F = \frac{\text{MST}_{\text{rt}}}{\text{MSE}}$$

- If  $H_0$  is true,  $F \sim F(\text{df}_1, \text{df}_2)$  (Central  $F$ )
- If a specific  $H_A$  is true  $F \sim F(\text{df}_1, \text{df}_2, \lambda)$  (Noncentral  $F$ )

# Degrees of freedom (CRD)

If we have  $t$  treatments and  $r$  reps each:

- $df_1 = t - 1$
- $df_2 = t(r - 1)$

## *i* Example 3.1: Skeleton ANOVA

SV	DF: 12 runners - 1 = 11 total
Shoe Type	$(4 - 1) = 3$
Runner(Shoe Type) → error	$(3 - 1)(4) = 8$

# The noncentrality parameter

$$\lambda = \frac{\sum_{i=1}^t (\mu_i - \bar{\mu})^2}{\sigma^2 / r}$$

- \_\_\_\_\_ = the mean of the  $i^{th}$  treatment group
- \_\_\_\_\_ = the overall mean
- \_\_\_\_\_ = the experimental error variance
- \_\_\_\_\_ = the number of replications of the  $i^{th}$  treatment group

## **Sanity check**

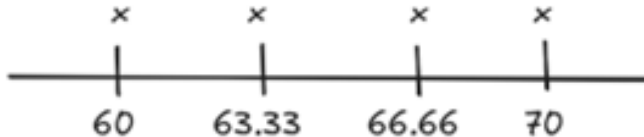
If all means are equal, then  $\lambda = 0$ .

# Example 3.1: Shoe Type

Recall, the things we need to calculate power are:

## 1. A specific alternative (e.g., $\delta$ )

An ultra-marathon running expert tells us that if there is a difference of 10 seconds between the runners lap times, then that would be of practical importance. Thus,  $\delta = 10$ .



## Example 3.1: Shoe Type

Recall, the things we need to calculate power are:

### **2. Replications per treatment ( $r$ )**

We are conducting a study with  $t = 4$  and  $r = 3$ .

# Example 3.1: Shoe Type

Recall, the things we need to calculate power are:

## 3. variability estimate ( $\sigma^2$ )

From an analysis of a pilot study (from Mod 2: CRD Notes), we found  $\hat{\sigma}^2 = 17.46$

### Pilot Studies

Pilot Studies can be super helpful for estimating experimental error variance ( $\hat{\sigma}^2$ ) and differences in group means – key pieces of a power analysis. Running a small version of your experiment gives you data to calculate the mean square error and get a sense of the effect size. This makes it easier to plan the full study and ensure your design has enough power to detect meaningful differences.

## Example 3.1: Shoe Type

---

Recall, the things we need to calculate power are:

### **4. Significance level ( $\alpha$ )**

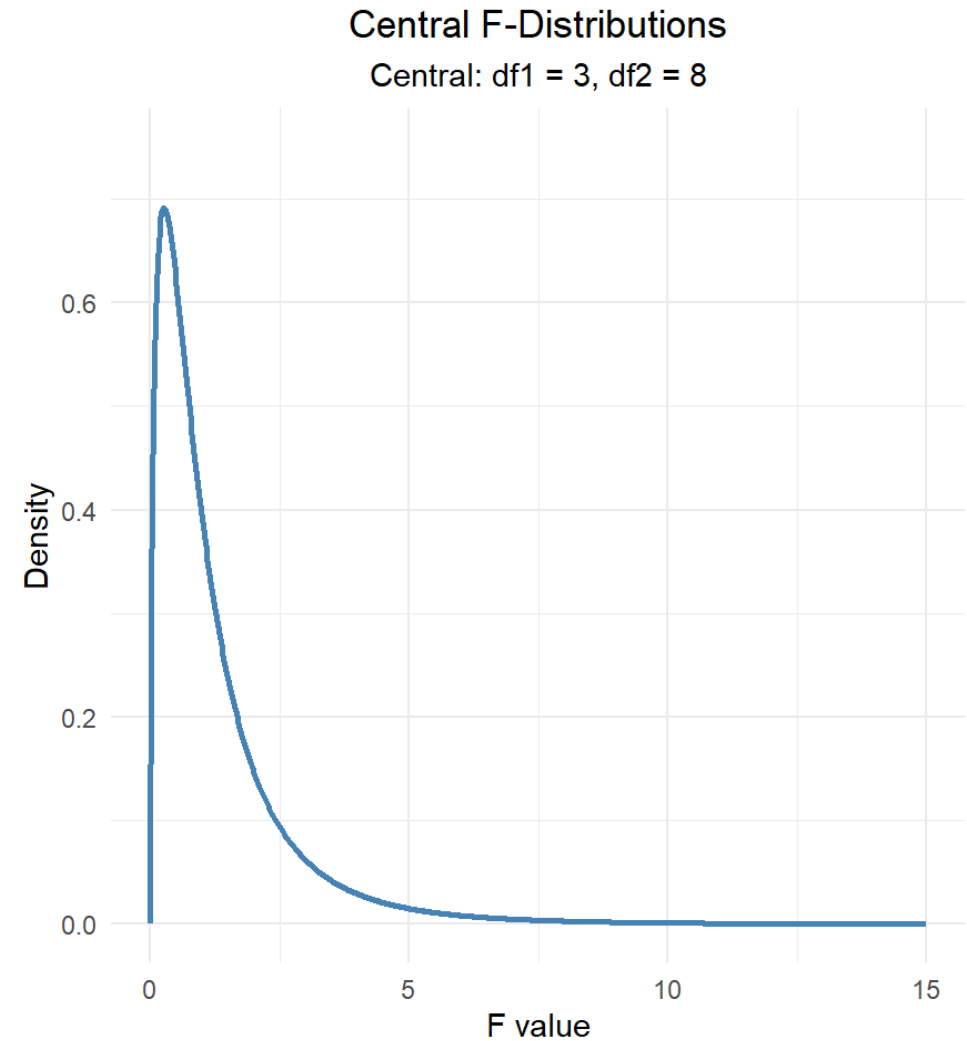
We get to set this at  $\alpha = 0.05$ .

# Power as area



# Power as an area

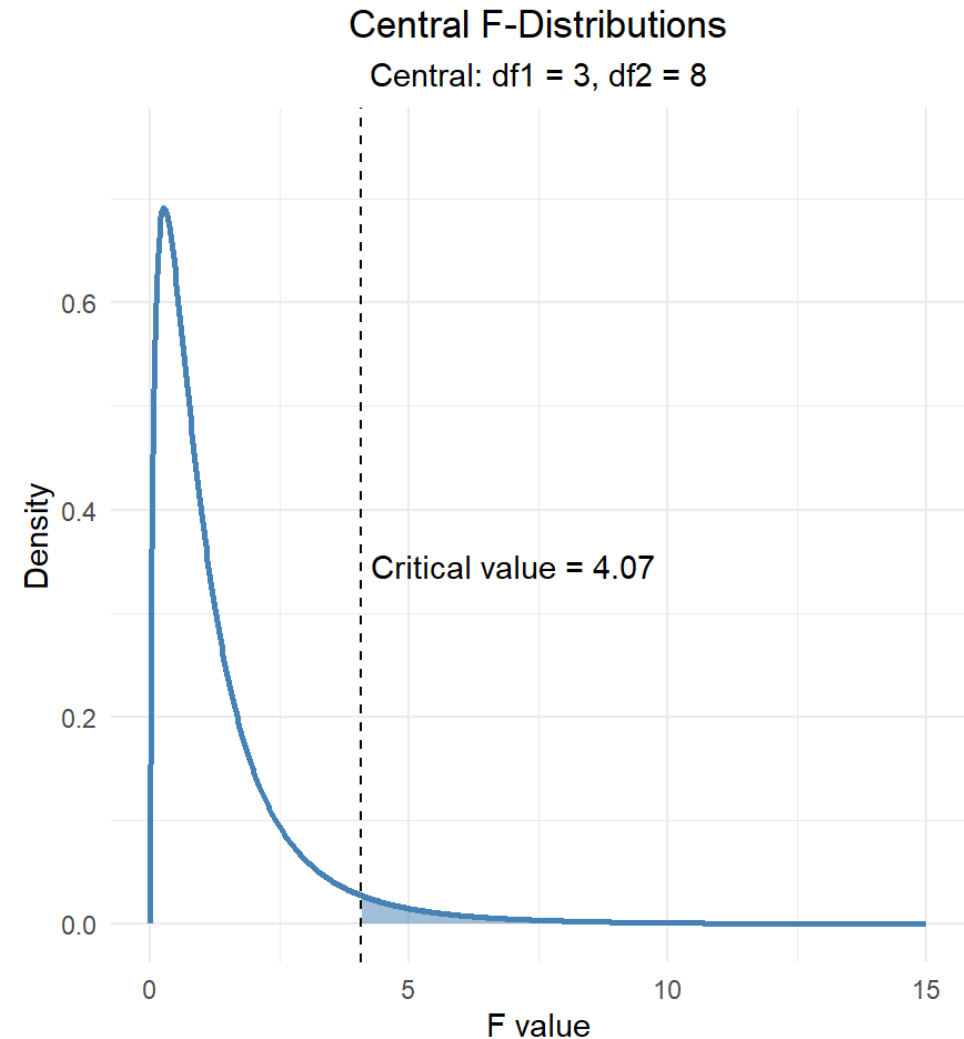
Consider the central F-distribution with these parameters:





# Power as area

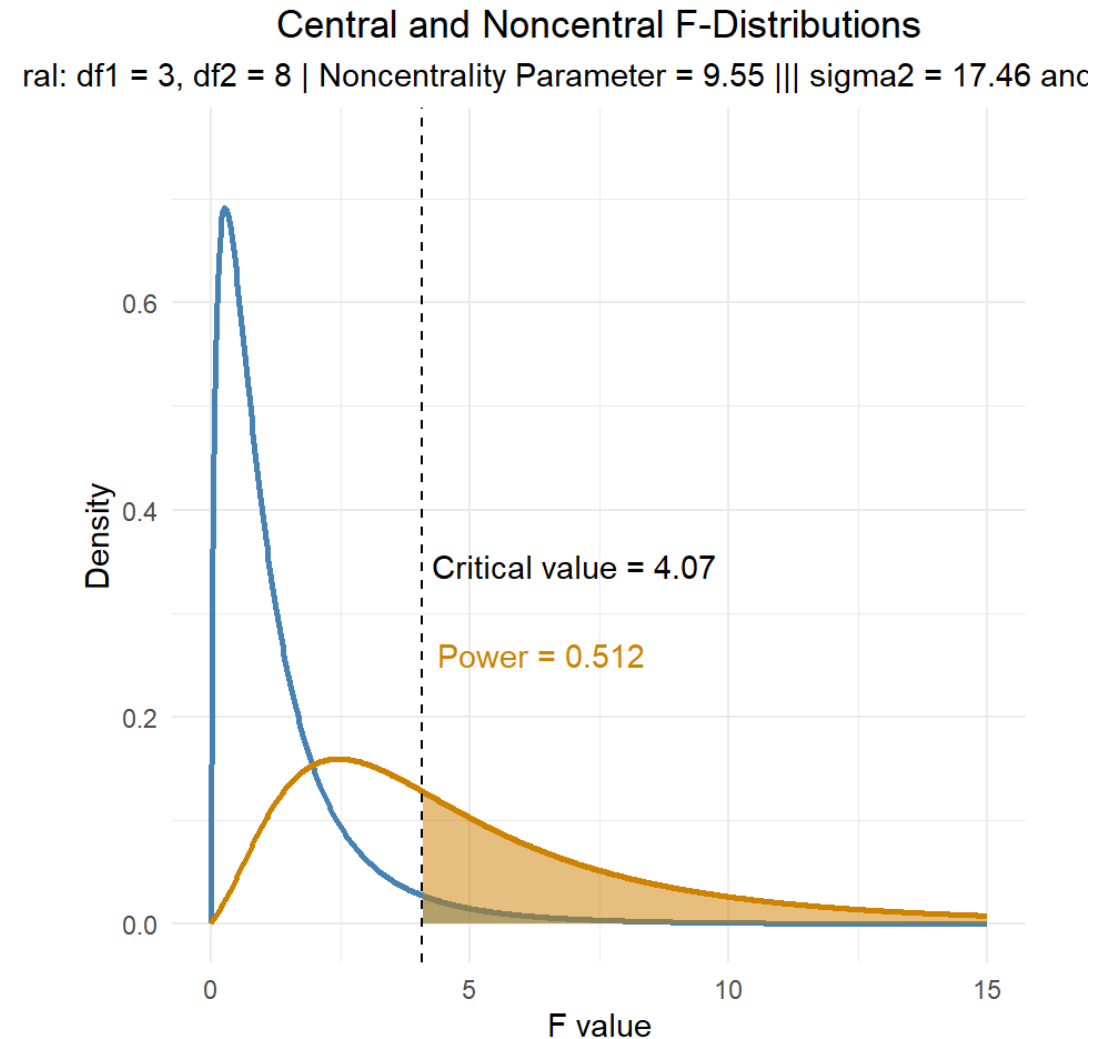
Recall, that the power represents the probability of **rejecting the null** when a specific alternative is true. So let's focus on the region where we'll reject the null hypothesis (the Rejection Region)





# Power as area

Now, we must incorporate the part about the **specific alternative being true**. If our particular alternative is true, then the F-statistic actually follows a non-central F-distribution with the aforementioned parameters. This is plotted on the same plot as the central F-distribution below:



The power for this test is about \_\_\_\_\_. That is, there is only a 51% chance of correctly rejecting the null hypothesis under these conditions. Do you think this is a good test?

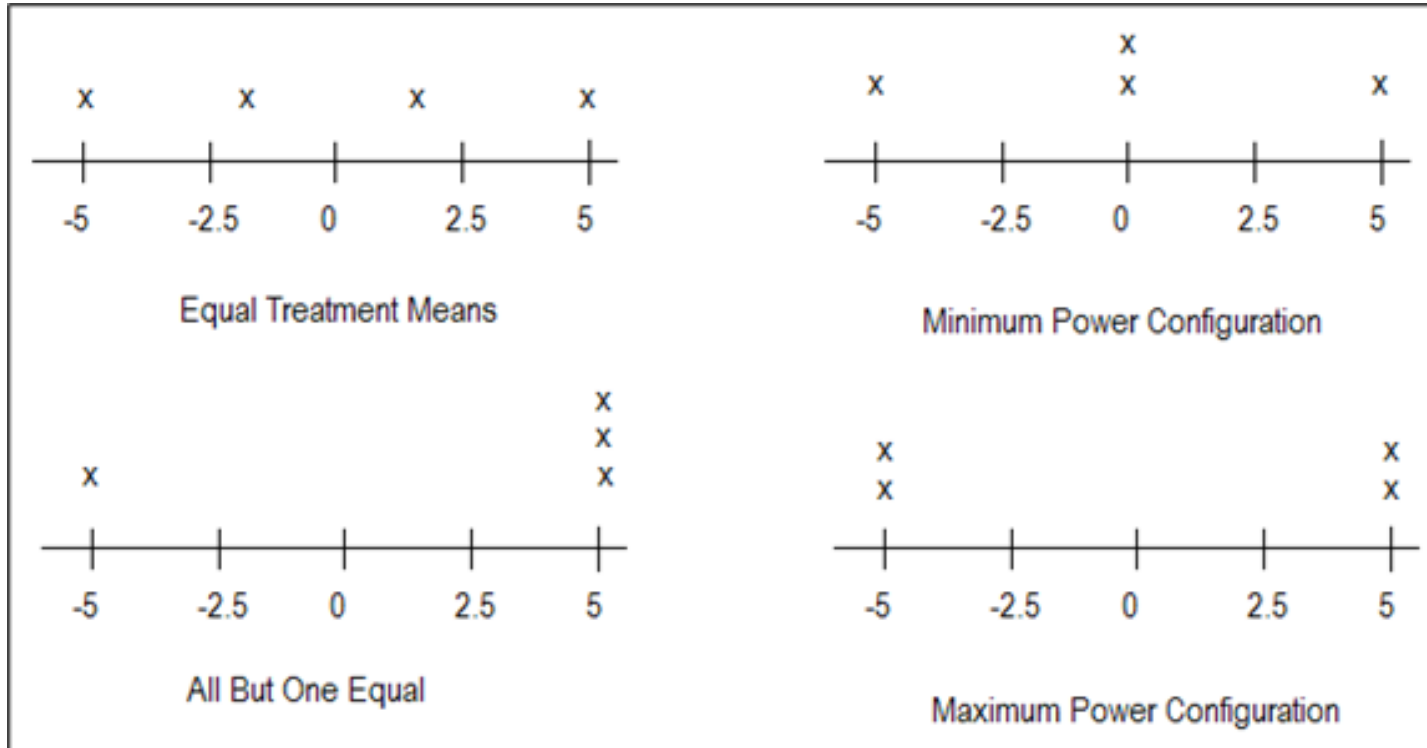
# Ideal Power

---

For most experiments, we would like the power of the test to be at least \_\_\_\_\_.

# Power Configurations

Even with the same max difference  $\delta$ ...





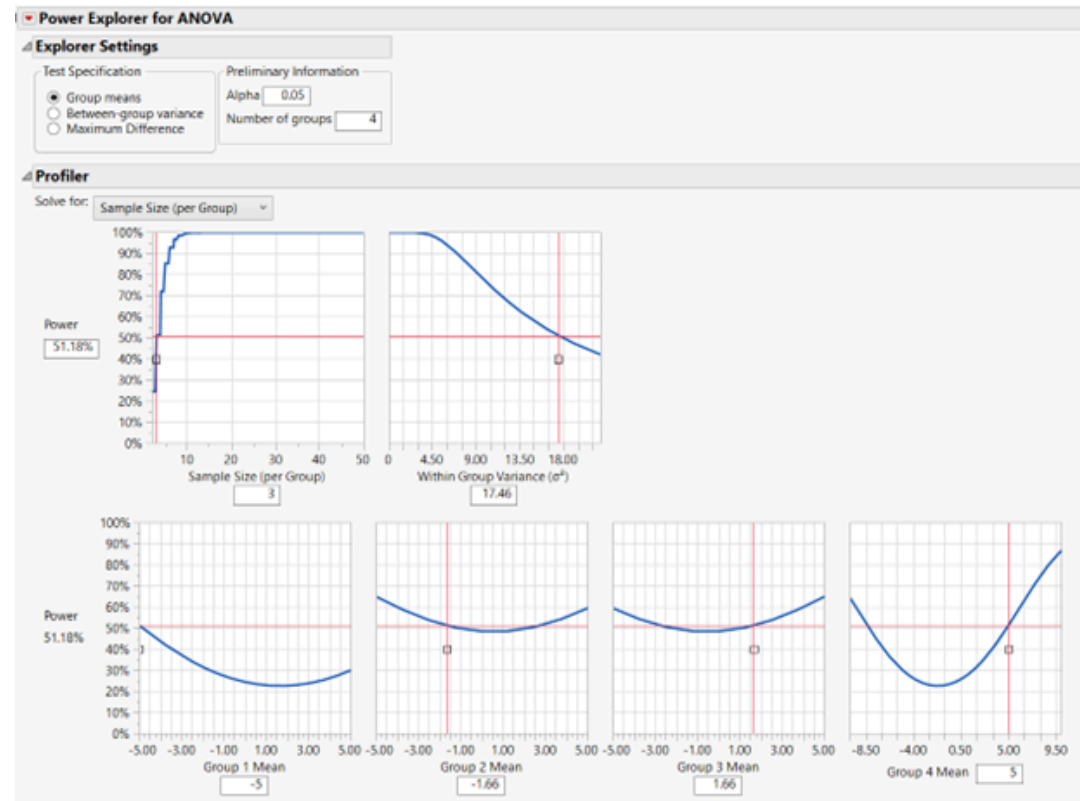
# Power depends on the configuration

Mean pattern	Power (approx.)
Maximum power	0.78
<b>Minimum power</b>	<b>0.47</b>
Equally spaced	0.51
All but one equal	0.65

# Using JMP

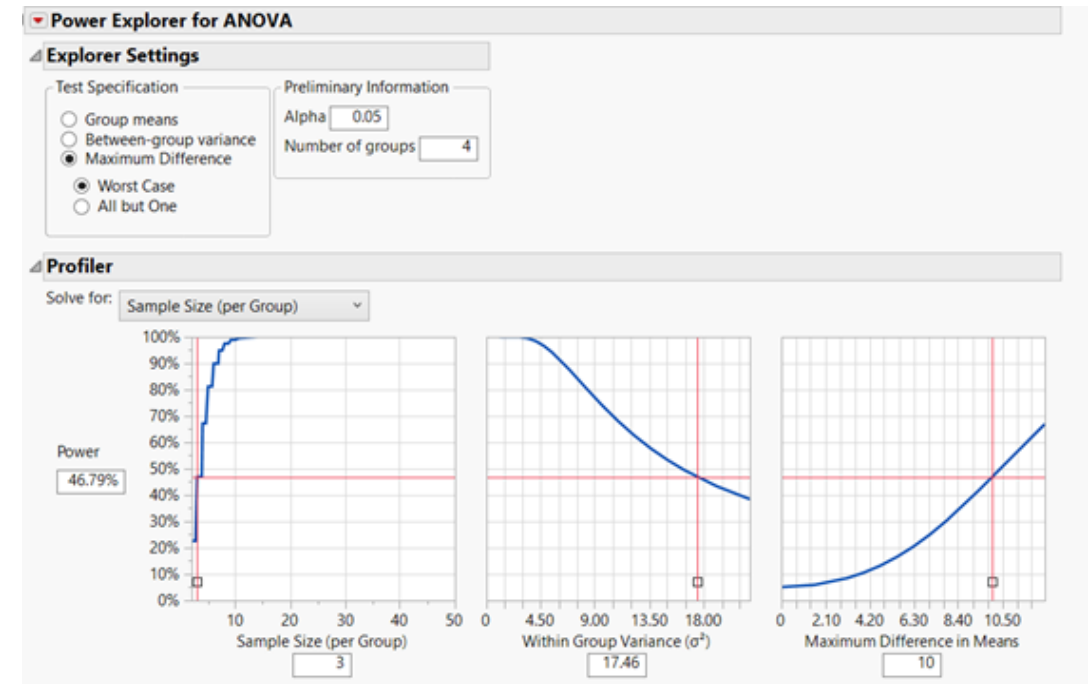
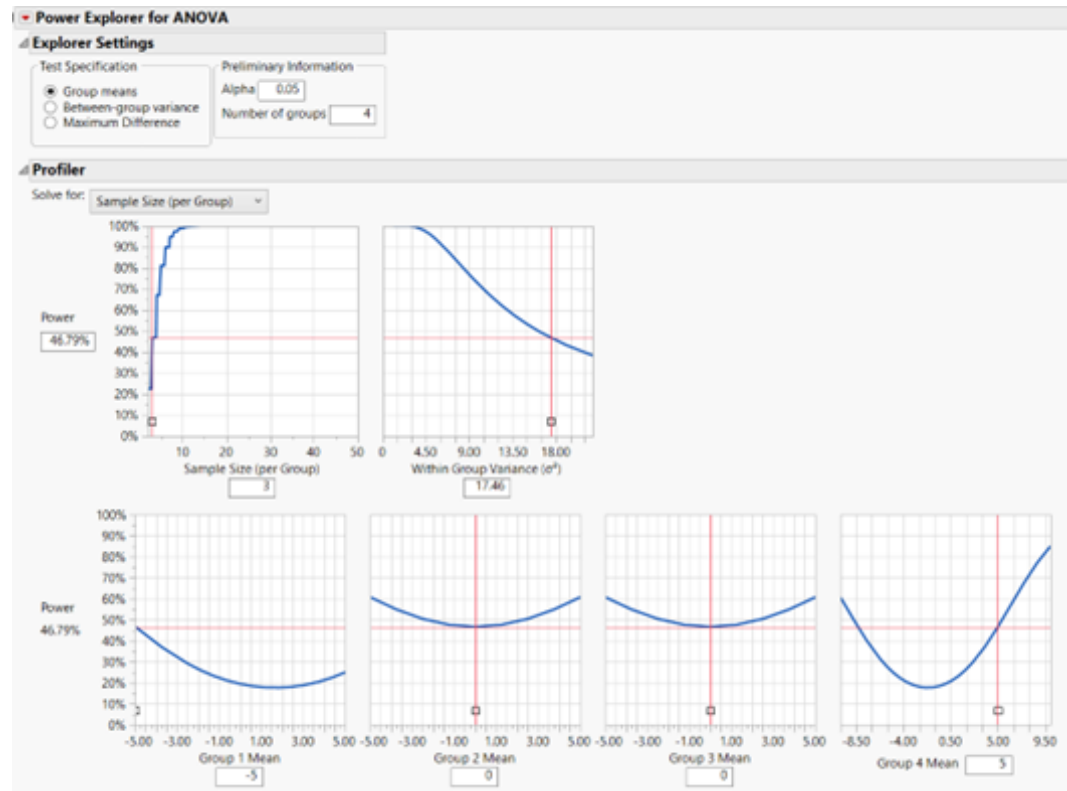
# JMP: Power Explorer

DOE > Sample Size Explorers > Power > Power for ANOVA



# JMP: Power Explorer

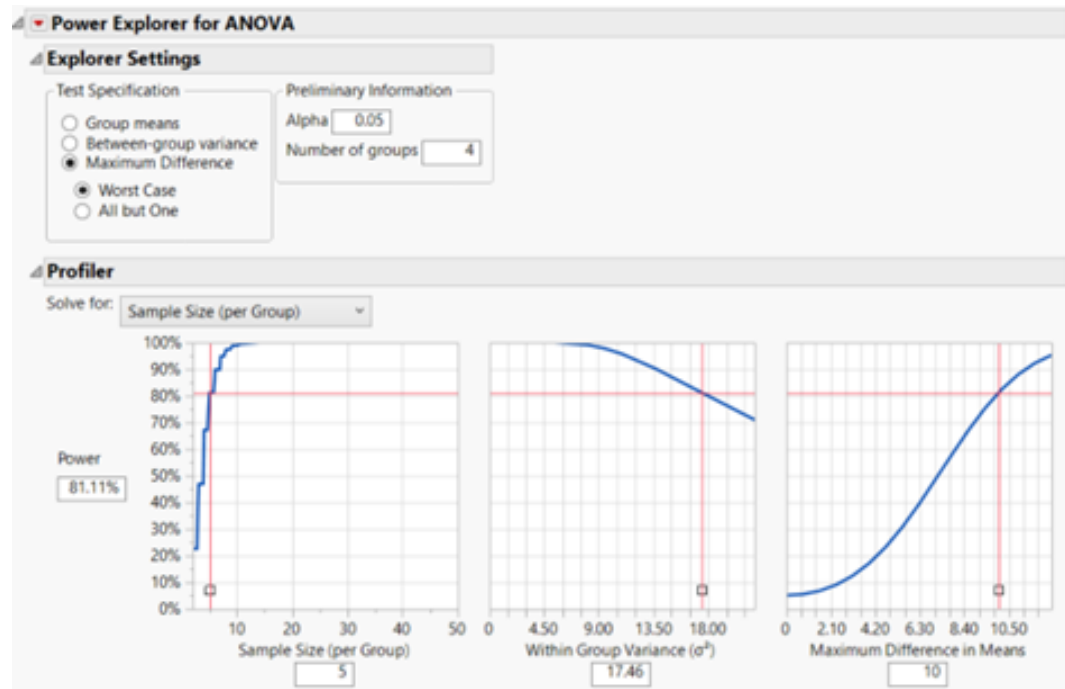
Minimum Power = Maximum Difference > Worst Case



# JMP: Power Explorer

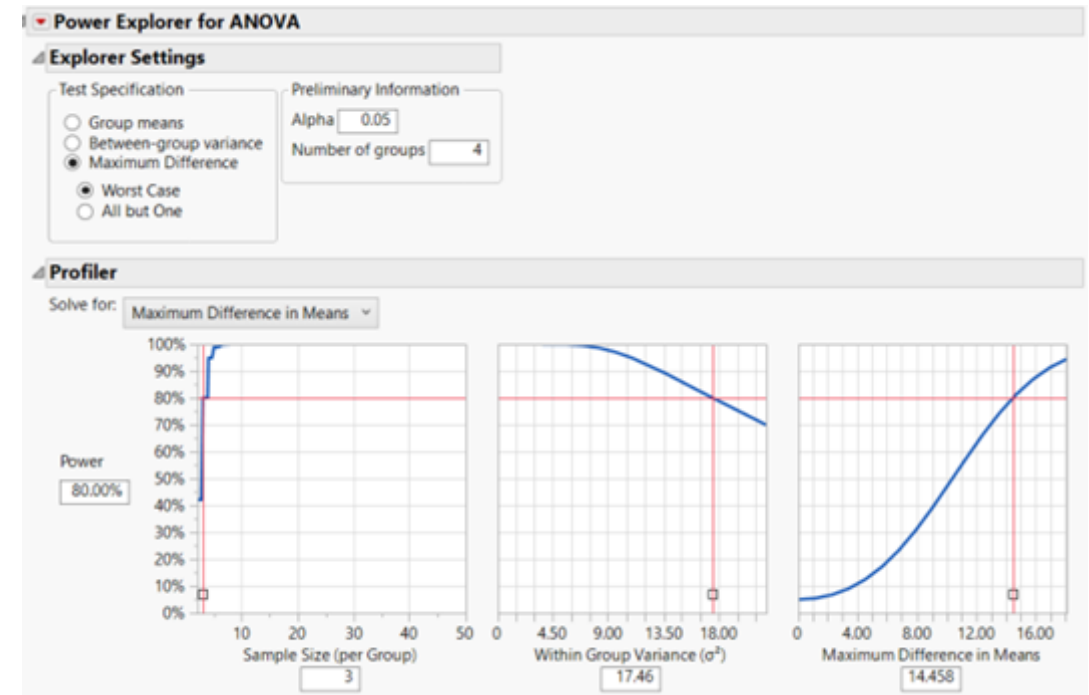
How many reps do we need?

Fix  $\alpha$ ,  $\sigma^2$ ,  $\delta$ , power  
Solve for  $r$



What difference can we detect?

Fix  $\alpha$ ,  $\sigma^2$ ,  $r$ , power  
Solve for  $\delta$



# Power in R

What is the the power of my test?

```
1 library(pwr2)
2 pwr.lway(
3   k = 4,           # Number of groups
4   n = 2,           # Sample size per group
5   alpha = 0.05,    # Significance level
6   delta = 10,      # Difference in group means
7   sigma = sqrt(17.46) # Standard deviation of obs
8 )
```

Balanced one-way analysis of variance power calculation

```
      k = 4
      n = 2
      delta = 10
      sigma = 4.178516
effect.size = 0.8461218
sig.level = 0.05
power = 0.2254578
```

NOTE: n is number in each group, total sample = 8 power = 0.225457781758005

What is the sample size I need?

```
1 ss.lway(
2   k = 4,           # Number of groups
3   alpha = 0.05,    # Significance level
4   beta = 1 - 0.80, # Type II error rate (1 - power)
5   delta = 10,      # Difference in group means
6   sigma = sqrt(17.46), # Standard deviation of obs
7   B = 100          # Number of iterations for numerical search
8 )
```

Balanced one-way analysis of variance sample size adjustment

```
      k = 4
sig.level = 0.05
power = 0.8
      n = 5
```

NOTE: n is number in each group, total sample = 20