# R in Action

## SECOND EDITION
*Data analysis and graphics with R*

ROBERT I. KABACOFF

# *Power analysis*

**This chapter covers**

- Determining sample size requirements
- Calculating effect sizes
- Assessing statistical power

As a statistical consultant, I'm often asked, "How many subjects do I need for my study?" Sometimes the question is phrased this way: "I have *x* number of people available for this study. Is the study worth doing?" Questions like these can be answered through *power analysis*, an important set of techniques in experimental design.

Power analysis allows you to determine the sample size required to detect an effect of a given size with a given degree of confidence. Conversely, it allows you to determine the probability of detecting an effect of a given size with a given level of confidence, under sample size constraints. If the probability is unacceptably low, you'd be wise to alter or abandon the experiment.

In this chapter, you'll learn how to conduct power analyses for a variety of statistical tests, including tests of proportions, t-tests, chi-square tests, balanced one-way ANOVA, tests of correlations, and linear models. Because power analysis applies to hypothesis testing situations, we'll start with a brief review of null hypothesis significance testing (NHST). Then we'll review conducting power analyses within R, focus-

ing primarily on the pwr package. Finally, we'll consider other approaches to power analysis available with R.

## 10.1   *A quick review of hypothesis testing*

To help you understand the steps in a power analysis, we'll briefly review statistical hypothesis testing in general. If you have a statistical background, feel free to skip to section 10.2.

In statistical hypothesis testing, you specify a hypothesis about a population parameter (your *null hypothesis,* or $H_0$). You then draw a sample from this population and calculate a statistic that's used to make inferences about the population parameter. Assuming that the null hypothesis is true, you calculate the probability of obtaining the observed sample statistic or one more extreme. If the probability is sufficiently small, you reject the null hypothesis in favor of its opposite (referred to as the *alternative* or *research hypothesis,* $H_1$).

An example will clarify the process. Say you're interested in evaluating the impact of cell phone use on driver reaction time. Your null hypothesis is Ho: $\mu_1 - \mu_2 = 0$, where $\mu_1$ is the mean response time for drivers using a cell phone and $\mu_2$ is the mean response time for drivers that are cell phone free (here, $\mu_1 - \mu_2$ is the population parameter of interest). If you reject this null hypothesis, you're left with the alternate or research hypothesis, namely $H_1$: $\mu_1 - \mu_2 \neq 0$. This is equivalent to $\mu_1 \neq \mu_2$, that the mean reaction times for the two conditions are not equal.

A sample of individuals is selected and randomly assigned to one of two conditions. In the first condition, participants react to a series of driving challenges in a simulator while talking on a cell phone. In the second condition, participants complete the same series of challenges but without a cell phone. Overall reaction time is assessed for each individual.

Based on the sample data, you can calculate the statistic $(\overline{X}_1 - \overline{X}_2)/(s/\sqrt{n})$, where $\overline{X}_1$ and $\overline{X}_2$ are the sample reaction time means in the two conditions, s is the pooled sample standard deviation, and n is the number of participants in each condition. If the null hypothesis is true and you can assume that reaction times are normally distributed, this sample statistic will follow a t distribution with $2n - 2$ degrees of freedom. Using this fact, you can calculate the probability of obtaining a sample statistic this large or larger. If the probability (p) is smaller than some predetermined cutoff (say $p < .05$), you reject the null hypothesis in favor of the alternate hypothesis. This predetermined cutoff (0.05) is called the *significance level* of the test.

Note that you use *sample* data to make an inference about the *population* it's drawn from. Your null hypothesis is that the mean reaction time of *all* drivers talking on cell phones isn't different from the mean reaction time of *all* drivers who aren't talking on cell phones, not just those drivers in your sample. The four possible outcomes from your decision are as follows:

- If the null hypothesis is false and the statistical test leads you to reject it, you've made a correct decision. You've correctly determined that reaction time is affected by cell phone use.

- If the null hypothesis is true and you don't reject it, again you've made a correct decision. Reaction time isn't affected by cell phone use.
- If the null hypothesis is true but you reject it, you've committed a Type I error. You've concluded that cell phone use affects reaction time when it doesn't.
- If the null hypothesis is false and you fail to reject it, you've committed a Type II error. Cell phone use affects reaction time, but you've failed to discern this.

Each of these outcomes is illustrated in the following table:

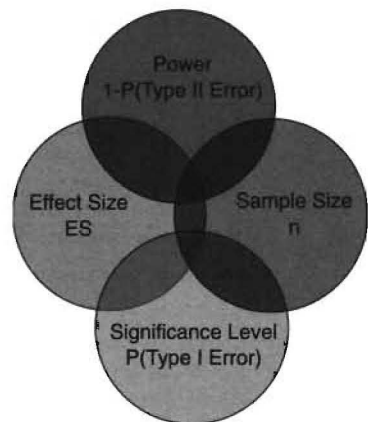| | | Decision | |
|---|---|---|---|
| | | **Reject $H_0$** | **Fail to Reject $H_0$** |
| **Actual** | **$H_0$ true** | Type I error | correct |
| | **$H_0$ false** | correct | Type II error |

### Controversy surrounding null hypothesis significance testing

Null hypothesis significance testing isn't without controversy; detractors have raised numerous concerns about the approach, particularly as practiced in the field of psychology. They point to a widespread misunderstanding of p values, reliance on statistical significance over practical significance, the fact that the null hypothesis is never exactly true and will always be rejected for sufficient sample sizes, and a number of logical inconsistencies in NHST practices.

An in-depth discussion of this topic is beyond the scope of this book. Interested readers are referred to Harlow, Mulaik, and Steiger (1997).

In planning research, the researcher typically pays special attention to four quantities (see figure 10.1):

- *Sample size* refers to the number of observations in each condition/group of the experimental design.
- The *significance level* (also referred to as *alpha*) is defined as the probability of making a Type I error. The significance level can also be thought of as the probability of finding an effect that is *not* there.
- *Power* is defined as one minus the probability of making a Type II error. Power can be thought of as the probability of finding an effect that *is* there.



**Figure 10.1  Four primary quantities considered in a study design power analysis. Given any three, you can calculate the fourth.**

- *Effect size* is the magnitude of the effect under the alternate or research hypothesis. The formula for effect size depends on the statistical methodology employed in the hypothesis testing.

Although the sample size and significance level are under the direct control of the researcher, power and effect size are affected more indirectly. For example, as you relax the significance level (in other words, make it easier to reject the null hypothesis), power increases. Similarly, increasing the sample size increases power.

Your research goal is typically to maximize the power of your statistical tests while maintaining an acceptable significance level and employing as small a sample size as possible. That is, you want to maximize the chances of finding a real effect and minimize the chances of finding an effect that isn't really there, while keeping study costs within reason.

The four quantities (sample size, significance level, power, and effect size) have an intimate relationship. *Given any three, you can determine the fourth.* You'll use this fact to carry out various power analyses throughout the remainder of the chapter. In the next section, we'll look at ways of implementing power analyses using the R package pwr. Later, we'll briefly look at some highly specialized power functions that are used in biology and genetics.

## 10.2  *Implementing power analysis with the pwr package*

The pwr package, developed by Stéphane Champely, implements power analysis as outlined by Cohen (1988). Some of the more important functions are listed in table 10.1. For each function, you can specify three of the four quantities (sample size, significance level, power, effect size), and the fourth will be calculated.

**Table 10.1  pwr package functions**

| Function | Power calculations for … |
|---|---|
| pwr.2p.test | Two proportions (equal n) |
| pwr.2p2n.test | Two proportions (unequal n) |
| pwr.anova.test | Balanced one-way ANOVA |
| pwr.chisq.test | Chi-square test |
| pwr.f2.test | General linear model |
| pwr.p.test | Proportion (one sample) |
| pwr.r.test | Correlation |
| pwr.t.test | t-tests (one sample, two samples, paired) |
| pwr.t2n.test | t-test (two samples with unequal n) |

Of the four quantities, effect size is often the most difficult to specify. Calculating effect size typically requires some experience with the measures involved and knowledge of past research. But what can you do if you have no clue what effect size to expect in a given study? We'll look at this difficult question in section 10.2.7. In the remainder of this section, we'll look at the application of pwr functions to common statistical tests. Before invoking these functions, be sure to install and load the pwr package.

## 10.2.1 t-tests

When the statistical test to be used is a t-test, the pwr.t.test() function provides a number of useful power analysis options. The format is

```
pwr.t.test(n=, d=, sig.level=, power=, type=, alternative=)
```

where

- n is the sample size.
- d is the effect size defined as the standardized mean difference.

$$d = \frac{\mu_1 - \mu_2}{\sigma} \quad \text{where} \quad \begin{aligned} &\mu_1 = \text{mean of group 1} \\ &\mu_2 = \text{mean of group 2} \\ &\sigma^2 = \text{common error variance} \end{aligned}$$

- sig.level is the significance level (0.05 is the default).
- power is the power level.
- type is a two-sample t-test ("two.sample"), a one-sample t-test ("one.sample"), or a dependent sample t-test ( "paired"). A two-sample test is the default.
- alternative indicates whether the statistical test is two-sided ("two.sided") or one-sided ("less" or "greater"). A two-sided test is the default.

Let's work through an example. Continuing the experiment from section 10.1 involving cell phone use and driving reaction time, assume that you'll be using a two-tailed independent sample t-test to compare the mean reaction time for participants in the cell phone condition with the mean reaction time for participants driving unencumbered.

Let's assume that you know from past experience that reaction time has a standard deviation of 1.25 seconds. Also suppose that a 1-second difference in reaction time is considered an important difference. You'd therefore like to conduct a study in which you're able to detect an effect size of d = 1/1.25 = 0.8 or larger. Additionally, you want to be 90% sure to detect such a difference if it exists, and 95% sure that you won't declare a difference to be significant when it's actually due to random variability. How many participants will you need in your study?

Entering this information in the pwr.t.test() function, you have the following:

```
> library(pwr)
> pwr.t.test(d=.8, sig.level=.05, power=.9, type="two.sample",
            alternative="two.sided")
```

```
Two-sample t test power calculation

            n = 34
            d = 0.8
    sig.level = 0.05
        power = 0.9
  alternative = two.sided
```

NOTE: n is number in *each* group

The results suggest that you need 34 participants in each group (for a total of 68 participants) in order to detect an effect size of 0.8 with 90% certainty and no more than a 5% chance of erroneously concluding that a difference exists when, in fact, it doesn't.

Let's alter the question. Assume that in comparing the two conditions you want to be able to detect a 0.5 standard deviation difference in population means. You want to limit the chances of falsely declaring the population means to be different to 1 out of 100. Additionally, you can only afford to include 40 participants in the study. What's the probability that you'll be able to detect a difference between the population means that's this large, given the constraints outlined?

Assuming that an equal number of participants will be placed in each condition, you have

```
> pwr.t.test(n=20, d=.5, sig.level=.01, type="two.sample",
             alternative="two.sided")

    Two-sample t test power calculation

            n = 20
            d = 0.5
    sig.level = 0.01
        power = 0.14
  alternative = two.sided
```

NOTE: n is number in *each* group

With 20 participants in each group, an a priori significance level of 0.01, and a dependent variable standard deviation of 1.25 seconds, you have less than a 14% chance of declaring a difference of 0.625 seconds or less significant ($d = 0.5 = 0.625/1.25$). Conversely, there's an 86% chance that you'll miss the effect that you're looking for. You may want to seriously rethink putting the time and effort into the study as it stands.

The previous examples assumed that there are equal sample sizes in the two groups. If the sample sizes for the two groups are unequal, the function

```
pwr.t2n.test(n1=, n2=, d=, sig.level=, power=, alternative=)
```

can be used. Here, n1 and n2 are the sample sizes, and the other parameters are the same as for pwer.t.test. Try varying the values input to the pwr.t2n.test function and see the effect on the output.