

# Power

Gary W. Oehlert

School of Statistics  
University of Minnesota

September 18, 2014

# Background

We have already defined type I and II errors.

Decision	Reality/State of nature	
	Null correct	Null false
Fail to reject	☺	Type II error
Reject	Type I error	☺

The type I error rate  $\mathcal{E}$  is easy to set, we just choose it.

Power is the probability of rejecting the null when the null is false.  
Power is the probability of declaring a difference when the difference is there (getting that lower right smiley).

Power is a much more difficult customer than  $\mathcal{E}$ .

You should design your experiments to have “appropriate” power.

- If the power is too low, then you're just wasting your time and resources running an experiment with no chance of finding what you are looking for.
- If the power is too high, then you are spending resources in this experiment that might be better spent somewhere else.

Appropriate power is probably in the .7 to .95 range, but it is situationally dependent.

Power for the F test comparing the separate means model with the single mean model depends on practically everything:

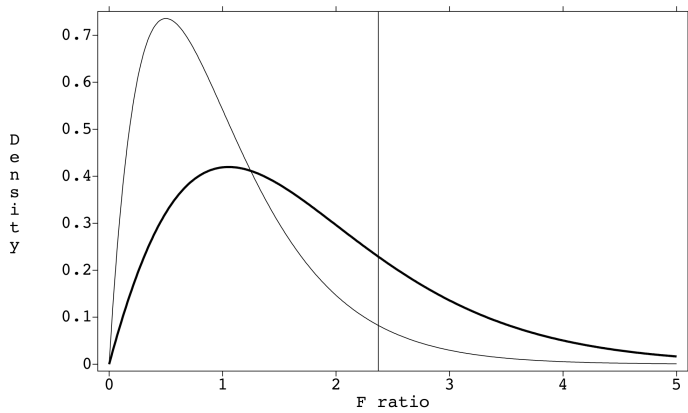
- The type I error rate  $\mathcal{E}$ .
- The numerator and denominator degrees of freedom for the F test; these obviously depend on  $g$  and  $N$ .
- The “non-centrality parameter  $\zeta$ , which itself depends on the sample sizes  $n_1, \dots, n_g$ , the non-null treatment means  $\mu_1, \dots, \mu_g$ , and the error variance  $\sigma^2$ .

Under the null, the  $F$  statistic follows a (central)  $F$  distribution with  $g-1$  and  $N-g$  df.

Combine this with  $\mathcal{E}$  and we get a critical value: reject for  $F$  statistics larger than the critical value. (Equivalently, any  $F$  in that range will have a  $p$ -value less than  $\mathcal{E}$ .)

When the null is false, the  $F$  statistic follows a non-central  $F$  distribution with  $g-1$  and  $N-g$  df. The distribution is shifted to the right, and  $\zeta$  controls the amount of shift to the right.

Probability of being to the right of the critical value is the power.



As you decrease  $\mathcal{E}$ , it becomes more difficult to reject the null (that moves the critical value to the right so you need a bigger F statistic to reject). For fixed  $g$ ,  $N$ , and  $\zeta$ , smaller  $\mathcal{E}$  leads to lower power.

$\zeta$  is a measure of how far the alternative state of nature is from the null. Replace the data with their respective means  $\mu_i$ , then fit the null model to these “data.” Get the residual SS and divide it by  $\sigma^2$ . That is  $\zeta$ .

- $\zeta$  increases if you increase the sample sizes.
- $\zeta$  increases if the error variance is smaller.
- $\zeta$  increases if the means  $\mu_i$  are further apart.

A formula for  $\zeta$ . Let  $\mu_j = \mu + \alpha_j$  where we use the  $\alpha_j$ s with  $\sum_j n_j \alpha_j = 0$ . Then

$$\zeta = \frac{\sum_{i=1}^g n_i \alpha_i^2}{\sigma^2}$$

Note: the expected value of the  $MS_{Trt}$  is

$$E[MS_{Trt}] = \sigma^2 + \frac{\sum_{i=1}^g n_i \alpha_i^2}{g - 1}$$

With  $E[MS_E] = \sigma^2$  you can start to see how these hook together with  $\zeta$ .



## Excuse me, but . . .

The discerning student will remark that  $\zeta$  depends on lots of stuff we don't know, like the  $\mu_i$ s and  $\sigma^2$ . *If we knew the  $\mu_i$ s, we wouldn't be doing the experiment in the first place!* So what gives?

In practice, power analysis and sample size selection are a big exercise in “Let's pretend” or “What if?”

We can control  $\mathcal{E}$ , and we can control  $n_1, \dots, n_g$ , but otherwise we are plugging in some hypothesized means and error variance and asking what the power would be for that state of nature.

To make power analysis useful, you must be able to specify some scientifically or practically meaningful set of alternative means  $\mu_i$ , and you must be able to make a guess as to how large the error variance is.

(Some people like to think of effects  $\alpha_i$  as multiples of  $\sigma$ , and while that is mathematically true, I think that is usually a cop out when specifying alternatives.)

Find alternative means where you can say, “If this were true, I would want to know about it,” and then design for those interesting alternatives.

Examples might be

- A doubling of the mutation rate is practically significant, so I want to design for that.
- An increase in MPG of 1 is relevant, so I will design for that.
- A 20% reduction in the serum concentration of a hormone is diagnostic, so I design for that.

Many (most?) granting agencies will require a power analysis before funding a proposal.

OK, but what about  $\sigma^2$ ? Some possibilities include:

- Variance from a pilot study.
- Variance from similar experiments in your lab or in the literature.
- Theoretical variances (possible for binomial counts and some other situations).
- Analytical variance of equipment (generally an underestimate of  $\sigma^2$ ).

It's probably best to do multiple power analyses that cover a range of plausible  $\sigma^2$  values.

Suppose that you have equal sample sizes  $n$ , and you think that any configuration of means where two means are  $D$  or more units apart is interesting.

The smallest value of  $\zeta$  for that description is

$$\zeta_0 = \frac{nD^2}{2\sigma^2}$$

Any  $\zeta$  for two means  $D$  units apart with sample sizes  $n$  will be at least as big as  $\zeta_0$ .

Thus the power for any of the other  $\zeta$ s will be at least as big as what you compute for  $\zeta_0$ .

# Sample size

You have chosen  $\mathcal{E}$ , you have some interesting values for the  $\mu_i$ s, and you have a pretty good idea what  $\sigma^2$  is.

Sample size analysis takes those and finds the smallest sample sizes  $n_i$  that will achieve a specified level of power.

In principle this involves computing power for a lot of different sample sizes and finding the one that is just big enough. In practice, we just use R.

# Confidence intervals

Another approach to sample sizes picks  $n$  so that confidence intervals are short enough.

For a contrast, we use the CI

$$\sum_{i=1}^g w_i \bar{y}_{i\bullet} \pm t_{\mathcal{E}/2, \nu} \sqrt{MS_E \sum_{i=1}^g \frac{w_i^2}{n_i}}$$

The margin of error is thus

$$MOE = t_{\mathcal{E}/2, \nu} \sqrt{MS_E \sum_{i=1}^g \frac{w_i^2}{n_i}}$$

where  $\nu$  is the df for MSE. The width of the interval is  $W = 2 \times MOE$ .

If we assume that the  $n_i$ s are all equal, we can solve to get:

$$n \approx \frac{t_{\mathcal{E}/2, \nu}^2 MS_E \sum_{i=1}^g w_i^2}{MOE^2}$$

We haven't done the experiment yet, so we don't know  $MS_E$ , and we will instead use a guess of  $\sigma^2$  as we did in power analysis.



We know our desired MOE, we know the  $w_i$ s, we have a guess for  $\sigma^2$  which we use as a guess for  $MS_E$ .

Compute  $n_0$  by substituting a normal percent point for the t-percent point.

$$n_0 \approx \frac{(\Phi^{-1}(1 - \mathcal{E}/2))^2 \sigma^2 \sum_{i=1}^g w_i^2}{MOE^2}$$

This gives you a starting point. Now start  $n$  at  $n_0$  and increment it until

$$n \geq \frac{t_{\mathcal{E}/2, g(n-1)}^2 \sigma^2 \sum_{i=1}^g w_i^2}{MOE^2}$$