

EPSE 592: Design & Analysis of Experiments

Ed Kroc

University of British Columbia

ed.kroc@ubc.ca

Statistical power

- The concept of *statistical power* is crucial for both designing a study and for interpreting a study that has already been conducted.
- *Power* is (informally) defined as the ability to detect non-zero effects (true positives)
- The *power*, or *sensitivity*, of a test is defined as

$$\Pr(p - \text{value} < \alpha \mid H_0 \text{ false}) = 1 - \beta,$$

where α is the *significance level* set by the researcher/journal and used to declare p-values “significant” or not under the traditional threshold approach.

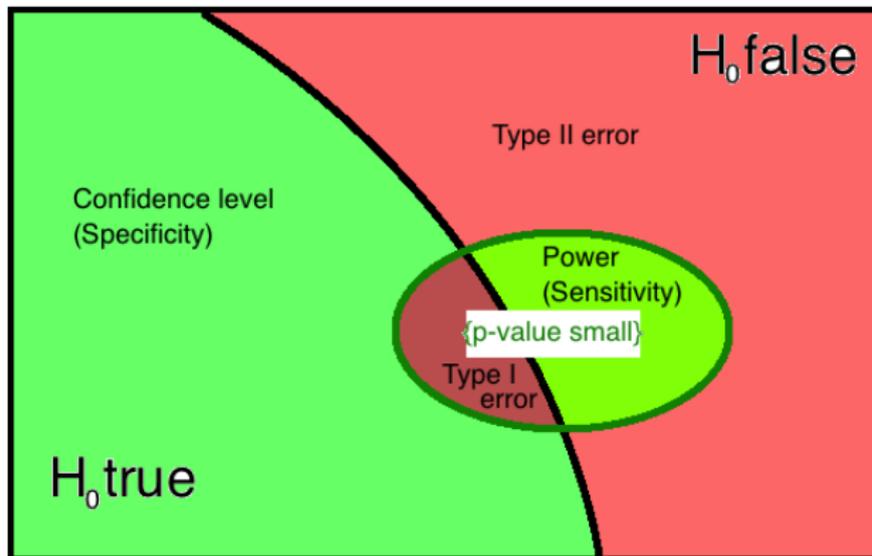
- Good studies will strive to have $1 - \beta \geq 0.80$. Most studies will have much lower power.

Statistical power

	H_0 true	H_0 false
data inconsistent with H_0	Type I error <i>false positive</i>	Correct decision <i>true positive</i>
data consistent with H_0	Correct decision <i>true negative</i>	Type II error <i>false negative</i>

	Given H_0 true	Given H_0 false
Pr(data inconsistent with H_0 ...)	α	$(1 - \beta)$
Pr(data consistent with H_0 ...)	$(1 - \alpha)$	β

Statistical power



- The only way to *simultaneously* decrease α and β (i.e. both kinds of errors) is to increase our sample size or choose a better (i.e. more powerful) statistical test.

Statistical power

- Statistical power is a function of many things:
 - Sample size (increasing sample size automatically increases power)
 - Population variability (less variation means more power)
 - Overall distribution of random phenomenon of interest (average effects in clustered or multi-modal distributions can be difficult to detect)
 - Type I error rate, α (increasing α automatically increases power)
 - *True, unobserved effect size* (bigger effect sizes are easier to find)
 - Type of statistical test/procedure used (e.g. nonparametric or robust procedures can be more powerful when data are non-normal)
 - Measurement error (noisier measurements produce more variability, so lead to less power)

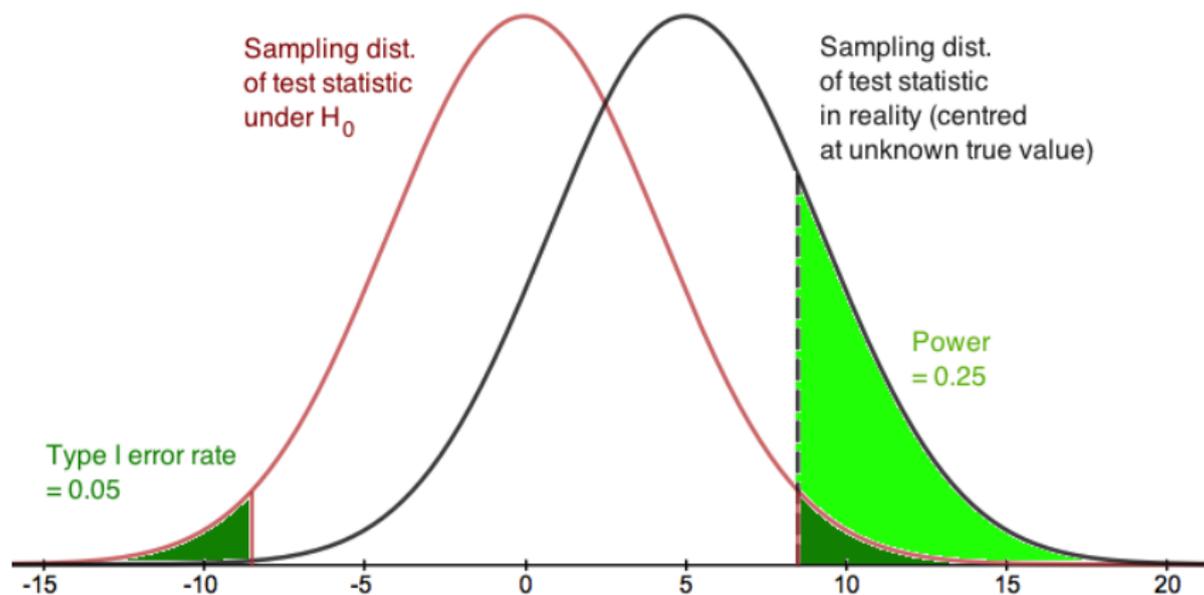
Statistical power

- When planning a study, power is considered to determine how large your *sample size* should be. This is called *power analysis* and generally proceeds as follows:
 - Identify the goal of the research study (e.g. testing if a new drug or intervention is more effective over current treatments)
 - Identify how you will measure the outcomes, effect size (e.g. mean difference between two treatment groups)
 - Use the previous literature to *reasonably estimate the variability* in your future study (e.g. similar drugs tested produced about a σ^2 variation in the response)
 - Decide on how you will analyze your outcomes (e.g. t-tests, ANOVAs, regression)
 - *Determine what effect size would be clinically important enough for you to care* (e.g. you want a new drug to be at least 20% more effective than current treatments)
 - Set your type I error rate α .
 - Set your desired power $1 - \beta$; i.e. your desired ability to detect the effect of clinical importance to you.

Statistical power

- Only after all this setup can we then estimate the necessary sample size to attain the desired power.
- This is a necessary step of virtually all medical research.
- This is often a necessary step to obtain funding for a proposed project. Why?
 - If you design a study that has a poor chance of detecting what you are trying to find, then why bother doing the study at all?
 - If your study has low power, but you end up finding a significant non-zero effect anyway, *it is likely that you are making a type I error.*
 - Moreover, if your study has low power but you end up finding a significant non-zero effect anyway, *your effect estimates are likely massively overinflated* (Type S and Type M errors).

Statistical power

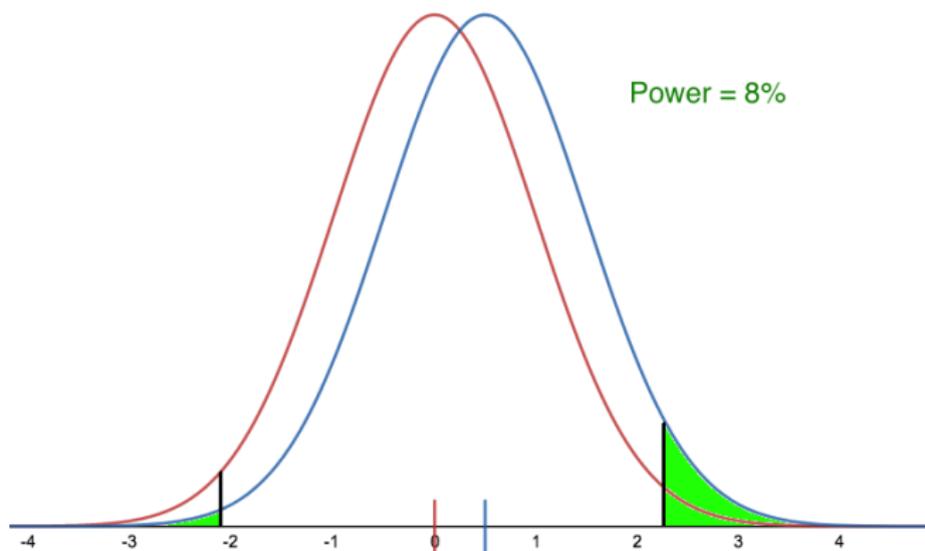


- Should *always* have this picture in mind when thinking about power.

Examples of study situations with different powers

Note: “small” and “large” are *relative* terms

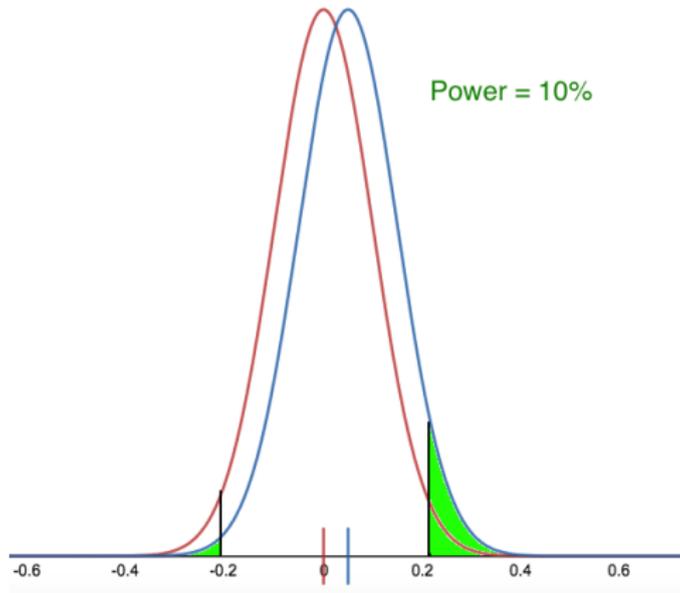
- Low power
- Small true effect size (0 vs. 0.5)
- Small sample size and/or large variance



Examples of study situations with different powers

Note: “small” and “large” are *relative* terms

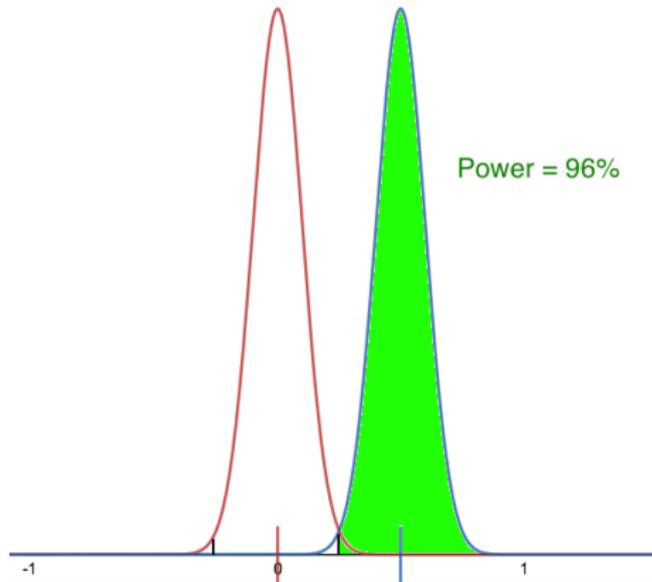
- Low power
- Very small true effect size (0 vs. 0.05)
- Large sample size and/or small variance



Examples of study situations with different powers

Note: “small” and “large” are *relative* terms

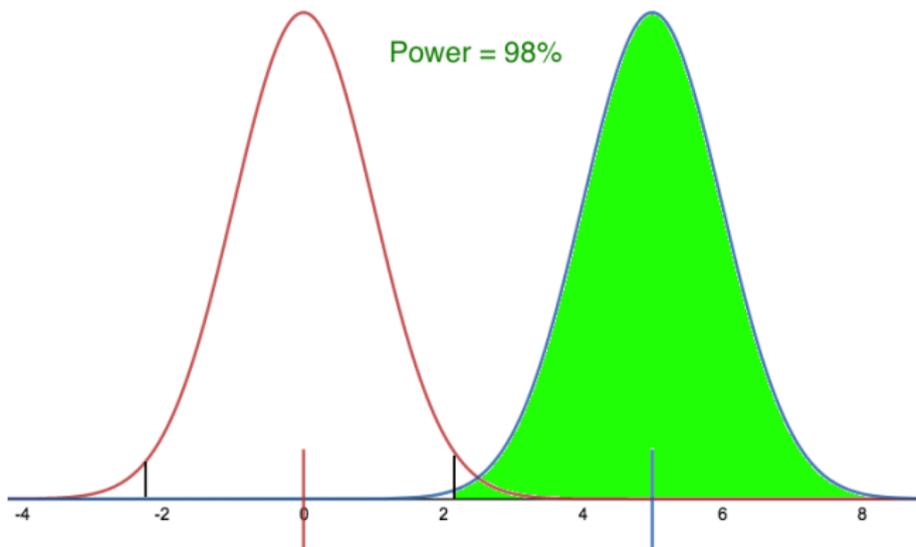
- High power
- Small true effect size (0 vs. 0.5)
- Large sample size and/or small variance



Examples of study situations with different powers

Note: "small" and "large" are *relative* terms

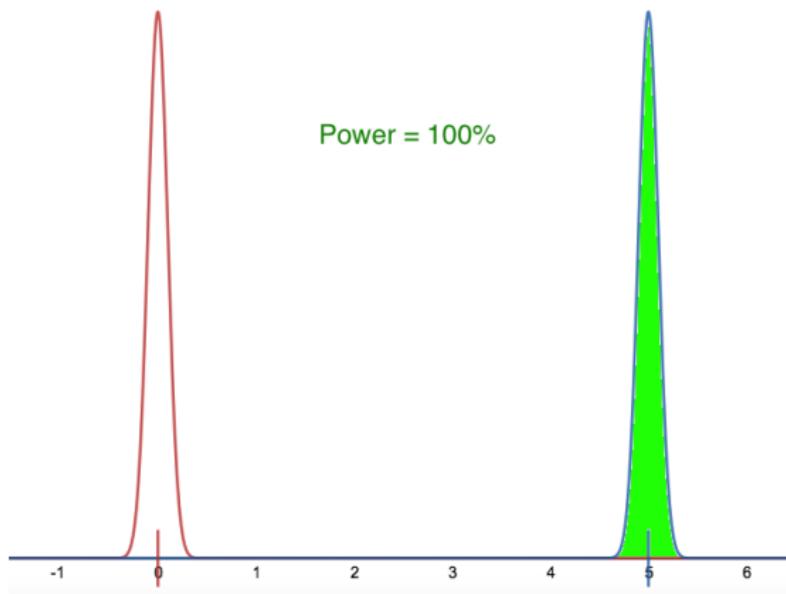
- High power
- Big true effect size (0 vs. 5)
- Small sample size and/or large variance



Examples of study situations with different powers

Note: “small” and “large” are *relative* terms

- Really high power (won't even require a statistical test of hypotheses)
- Big true effect size (0 vs. 5)
- Large sample size and/or small variance



How to calculate statistical power

- For simple scenarios, power can be calculated analytically (i.e. by hand). **But we rarely study simple scenarios.**
- *Lots* of software exists that claims to calculate power for you (e.g. SPSS, G*Power); but all of it relies on *the simple scenarios that rarely apply in practice.*
- In particular, software nearly always relies on an assumption of *perfectly normal data.*
- Practically, this means that sample size estimates can be distorted (very, very bad!)
- Usually no software or analytical options available for complicated study designs.
- What to do?

How to calculate statistical power

- What to do? **Must simulate (i.e. perform a simulation study) to perform power analysis.**
- Simulation allows you to tailor a sample size estimate to the exact specifics of any study design.
- Simulation requires semi-decent programming capabilities.
- If you don't have these skills, *seek a statistician's help!*

Effects of low power on interpretation of analytical output

Low power can come from many different sources. In practice, the three most common are:

- Small sample sizes (overall, or within groups).
- Large variability (overall, or within groups, or due to noisy measurements).
- Small *true* effect sizes.

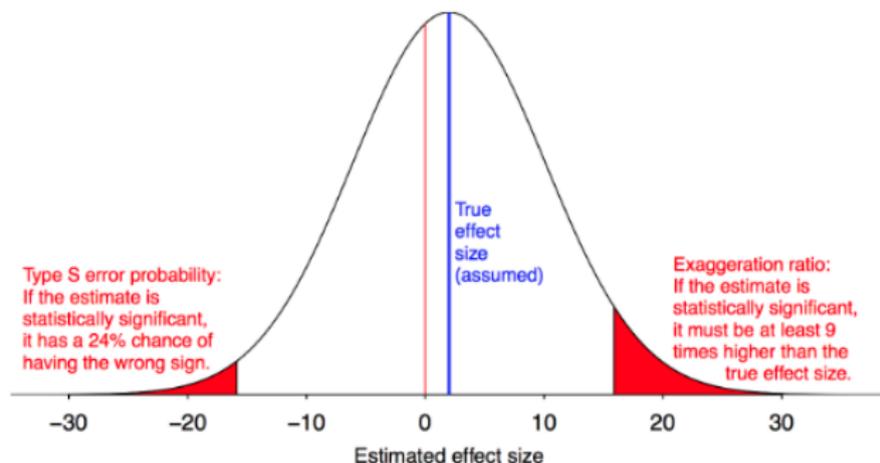
The first two sources are easy to see. The last (small true effect sizes) is difficult and somewhat subjective to assess, but absolutely crucial.

Effects of low power on interpretation of analytical output

True effect sizes are *unobserved*, but crucial to interpretation:

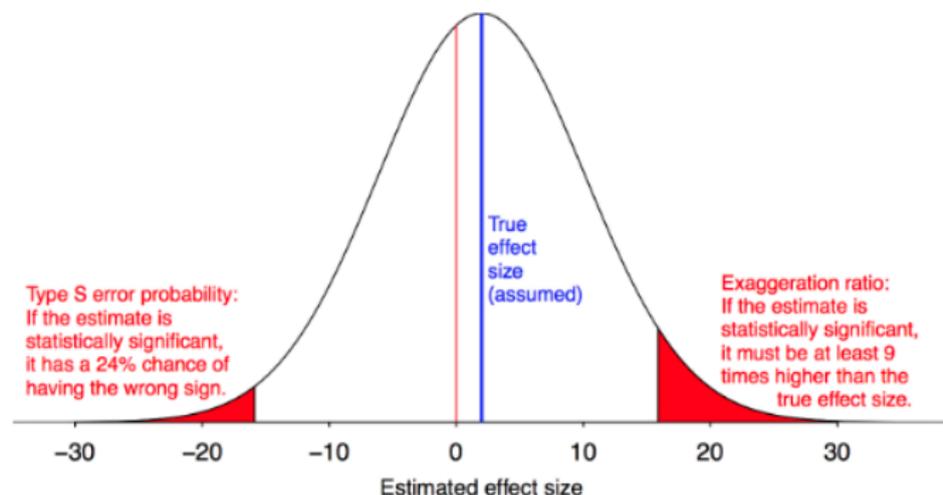
- We never actually know the *true* effect size (if we did, we wouldn't have to perform a study to estimate it).
- A plausible true effect size depends on the *prior believability of a particular alternative hypothesis*.
- In social science, many of our effects of interest will be small, *especially when compared to the effects of other variables of little or no interest*.
- *Evaluating the power of a study retrospectively requires an informed assessment of how plausible you would find certain effect sizes*.
- **Note:** some applied practitioners and software (e.g. SPSS) will talk about “retrospective power” or “post hoc power analysis”; they do *not* mean what we are talking about (usually, they mean gibberish).

Effects of low power on interpretation of analytical output



- This is a graphical representation of a t-test comparison of means.
- The *statistical power* here is 6%.
- In this example, true effect size (marked by blue line) is very small.
- Red regions represent values for “significant” test statistics (and so, p-values)

Effects of low power on interpretation of analytical output

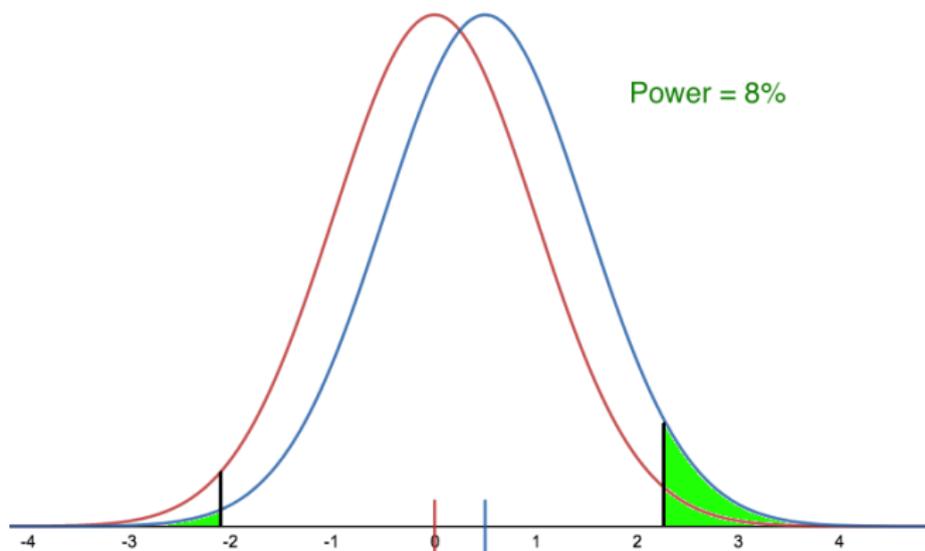


- But then finding a significant result would mean:
 - the estimated effect size is at least 9 times too big (Type M error)!
 - the estimated effect size has the wrong sign about 25% of the time (Type S error)! [See Gelman & Carlin (2014) for more info.]

Examples of study situations with different powers

Note: “small” and “large” are *relative* terms

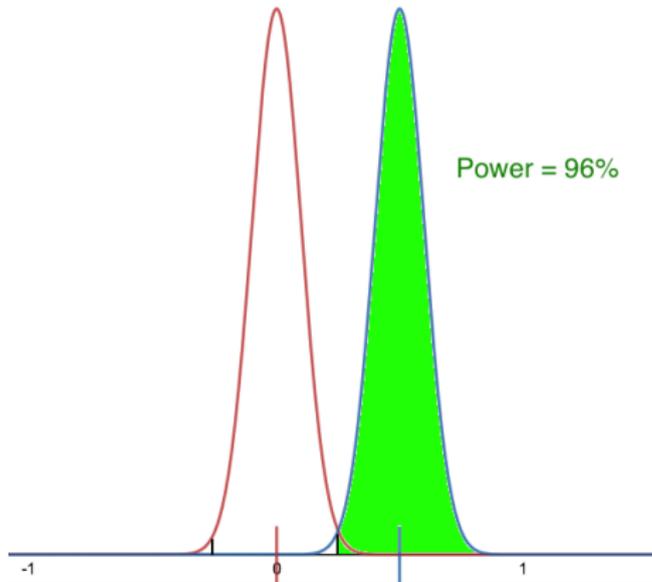
- Low power = bad estimates if significant
- Small true effect size (0 vs. 0.5)
- Small sample size and/or large variance



Examples of study situations with different powers

Note: “small” and “large” are *relative* terms

- High power = good estimates if significant
- Small true effect size (0 vs. 0.5)
- Large sample size and/or small variance



Effects of low power on interpretation of analytical output

In low-powered studies:

- Significant results are often meaningless.
- Significant results *will* yield estimates that are wildly inaccurate.
- Seemingly small things like measurement error, sampling variability, or minor experimental imperfections become magnified.
- Results are often entirely driven by statistical “noise”.

- Case study: Durante et al. 2013