

Some final words about hypothesis testing

Confidence intervals and hypothesis tests: two sides of the same coin?

In hypothesis testing, we construct a null distribution.

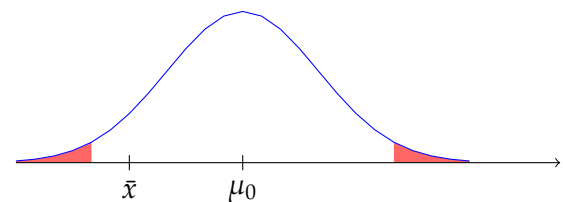
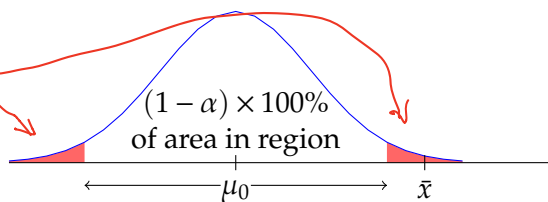
randomization
test.

$$H_0: \mu = \overset{\text{proposed number}}{\mu_0}$$

$$H_a: \mu \neq \mu_0$$

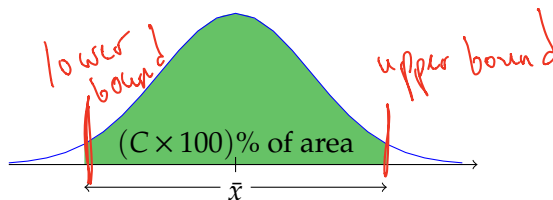
- often it has appeared to be symmetric, bell-shaped (normal?).
- null value is the mean/center
- in setting α we fix the area of the **rejection region** (two tails, colored red)
- Our test statistic may be in a tail (~~rejection~~ rejection region), or in the nonrejection region.

red region
has area
 $= \alpha$

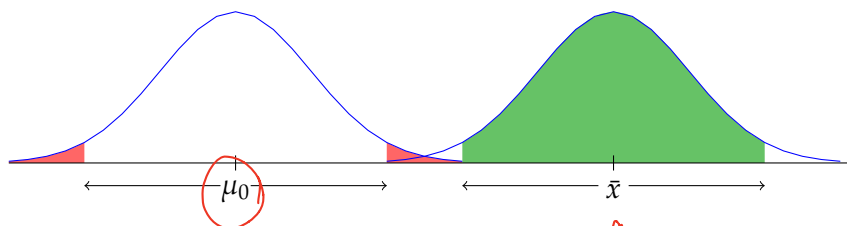


In confidence interval construction (centered interval method), we

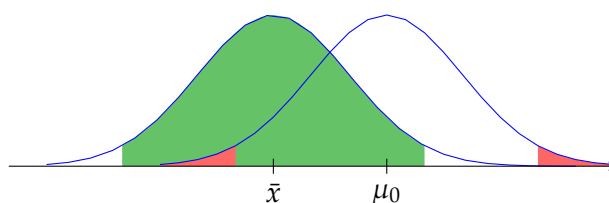
- use the point estimate as the center of a region (green).
- use the confidence level C to decide what portion to include in the region.



Now, in theory, whenever $C = 1 - \alpha$, the width of the nonrejection (uncolored) region of the null distribution is approximately the same as the width of the confidence interval. In this case, the composite picture might be like this



or like this



sample statistic

So, say a random sample is collected. The sample statistic, perhaps \bar{x} , could be used to construct a confidence interval for the parameter, perhaps μ , but it also could be used as a test statistic in the 2-sided test of hypotheses

$$H_0: \mu = \mu_0, \quad H_a: \mu \neq \mu_0.$$

But no matter which of these it is, confidence interval or hypothesis test, the one informs the other. Examples of the ways include these:

- If μ_0 is not in a 95% confidence interval, then the P -value of the hypothesis test is smaller than 0.05.
- If μ_0 is in a 90% confidence interval, then the P -value of the hypothesis test is larger than 0.1.
- If the P -value from the hypothesis test is 0.07, then μ_0 is in the 95% confidence interval, but not in the 90% confidence interval.

In case where C (conf. level) and α (significance level) add to 100%, get simultaneously

- sample statistic leads you to not reject H_0 along w/ null value lying the CI
- sample statistic is in rejection region (P -value smaller than α) along w/ null value lying outside CI.

Say we build a 95% CI for $\mu_1 - \mu_2$ 95% CI
[0.174, 1.562]

Get SE from bootstrapping say

$$\text{lower bound} = (\bar{x}_1 - \bar{x}_2) - 2SE = 0.174$$

$$\text{upper bound} = (\bar{x}_1 - \bar{x}_2) + 2SE = 1.562$$

Q: What conclusions can you make about Hypotheses

$$H_0: \mu_1 - \mu_2 = 0 \quad \text{vs.} \quad H_a: \mu_1 - \mu_2 \neq 0?$$

Cautions about multiple testing

Remember what was said earlier: setting $\alpha = 0.05$ for all your hypothesis tests means that, in cases where the null hypothesis is true, you will commit Type I error 5% of the time, *mistakenly* rejecting H_0 . That's a Type I error rate of 1-in-20. Any researcher conducting numerous statistical tests with $\alpha = 0.05$ should keep this in mind, and should maintain a healthy suspicion if about 5 percent of her tests are yielding statistically significant results. If, over the last year, 40 hypothesis tests have been conducted and 3 have been statistically significant, that is right near what we might expect to happen even if *none* of the null hypotheses in those tests of significance have been false.

$$\frac{3}{40} = 0.075 \quad \text{or} \quad 7.5\% \text{ of your results / experiments were stat. sign.}$$

Statistical significance is different from practical importance

Referring to the picture above, statistical significance amounts to our test statistic being far enough from the null value (μ_0) that it lands in the rejection region, nothing more. This may be evidence enough to reject the null hypothesis in favor of the alternative $H_a: \mu \neq \mu_0$, but it does not necessarily follow that the true value of μ is far away. You may have evidence that is statistically significant for showing that a drug does not leave blood pressure unchanged in those who suffer from high blood pressure even if its effect is only to decrease systolic pressure by 1.

Ellenberg has useful illustrations of hypothesis testing in Chapters 6 and 7 of his book, "How Not to Be Wrong: The Power of Mathematical Thinking." While I mention here his point that the word *significance* in statistics is meant in a technical sense that English language speakers are likely to misconstrue, I will let Ellenberg have the burden of hammering it home.

Many research questions

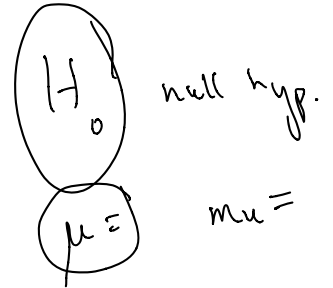
- involve two variables
- ask (ultimately) is there an association between these variables
- many get boiled down to questions about parameters

discuss null/alt. hypothesis in situations where

explanatory is binary categorical, so is response

explanatory is binary categorical, response is quantitative

explanatory is quantitative, so is response



- practice for 2-variable situations
 - 1 higher family income goes with higher grades (gpa) in school
 - 2 starting salaries for men are higher than starting salaries for women
 - 3 mean BMI is higher in America than in Europe
 - 4 rate of births is different on weekend days than weekdays
 - 5 sales resulting from radio ads are higher than for TV ads
 - 6 men who take aspirin have lower instance of heart attack
 - 7 women are in favor of gun control at a different level than men
 - 8 "sleep on it" results better in word-memorization/retention than caffeine
- coming later: bivariate situations when explanatory is categorical nonbinary
- univariate situations

researcher may question some widely-held(?) belief

examples:

 - ESP test (George has ESP vs. no he doesn't)
 - college students sleep less than 7 hours per night
 - is there more than 100 ppb of BPA in cans of tomatos?