

# Video 4: p-hacking

wbg231

December 2022

## introduction

- video link
- p values are everywhere in science
- they are often a request for publication because they show that what we see is not due to random fluctuations
- but they are not the only metric that should be taken into account
- a small p value neither implies practical significance nor causal effects

## practical significance

- we want to evaluate the cure rate of two really expensive drugs with a lot of side effects.
- in group 1A, where people do not receive the drug 30 out of 100 people recover
- in group 2A, where people receive drug A 52 out of 100 recover
- in group 1B, where people do not receive the drug 30,000 out of 100,000 people recover
- in group 2B, where people receive drug A 30650 out of 100,000 recover
- this seems to suggest that drug two is a lot less effective than drug 1.
- let's apply a two sample z test
- null there is no difference between the control and treatment group we assume all data are iid bern with cure rate parameter  $\theta$
- test statistic difference in cure rate between treatment and control groups
- under null the test stat  $\tilde{t} \sim \mathcal{N}(0, \theta(1 - \theta)(\frac{1}{n_{\text{treatment}}} + \frac{1}{n_{\text{control}}}))$

- notice that our variance depends a lot on the number of people in the study
- for drug one test stat is  $t_{data} = .2$  and this yields a very small p-value
- for drug two our test stat is 0.007 but with a much smaller test stat, so we have the same p value in the first trial
- this is interesting as they have the same significance level despite the fact that drug one makes much more of a difference in actual cure rate

### what does this mean

- both results equally unlikely under the null
- both increase the cure rate
- but they increase the cure rate by different amount
- so how do we quantify this difference?
- we can do a confidence interval for the difference in cure rate

### confidence interval

- let the true control cure rate be  $\theta_c$
- number of cured control subjects  $\tilde{k}_c$
- this is the number of cured patients in the control group is distributed as a binomial with parameters  $n_C, \theta_c$
- we can apply a gaussian approximation to the binomial with mean  $n_C\theta_c$  and variance  $n_C(\theta_c)(1 - \theta_c)$
- thus our observed control cure rate  $\frac{\tilde{k}_c}{n_c} \sim \mathcal{N}(\theta_c, \frac{\theta_c(1-\theta_c)}{n_c})$
- we can do the same thing for the treatment rate
- so we can think of the difference between the cure rate in the treatment and cure rate as distribution as  $\sim \mathcal{N}((\theta_t - \theta_c), \frac{\theta_c(1-\theta_c)}{n_c} + \frac{\theta_t(1-\theta_t)}{n_T})$
- so we can build a gaussian confidence interval around a gaussian rv  $\tilde{a}$  using mean  $\mu, \sigma^2$
- as  $\tilde{I}_{1-\alpha} := [\tilde{a} - c_\alpha\sigma, \tilde{a} + c_\alpha\sigma]$  where  $c_\alpha$  is some constant
- we see that the ci for drug 1 is between 8 and 35 percent
- drug 2 on the other hand as a difference between .25 and 1.05%
- so these two have the same p-value, but only drug 1 has a real effect

### **statistical vs practical significance**

- so more or less, if there is a test with a tone of power it can pick an effect that is so small that it likely does not matter in reality

### **covid example**

- an over powered study can make almost any effect size significant so it is important to think if the difference is practically significant

### **publication bias**

- we would expect the null to hold. ie that pizza does not cure covid-19
- so we will expect to see p-value uniformly distributed between 0 and 1
- so we would expect 5% of studies to find a false positive
- but if we do a lot of tests say 100 studies, then we would still expect 5%
- we would expect to see 5 false positives, and 95 true negatives
- this is not a big deal if all of these results were published
- but it is much easier to publish a false positive than a true negative
- so when this is done intentionally this is called p-hacking
- part of the issue with this is that it may not be reproducible later
- so this is a reproducibility issue.
- there are a lot of ways to fudge your stats in order to get statistically significant (this is unethical)
- it is an incentive problem

### **does p hacking happen in practice**

- the short answer is yes
- if we look at the distribution of p values on pub med,
- there is evidence that there are more studies published with lower p values