


# Day 03

## P-hacking Publication Bias Multiple testing

- Hypothesis testing
- P-value
- P-hacking, Publication Bias
- Multiple hypothesis testing

# Statistical hypothesis testing

## Abstract

Formula display: ☒ MathJax 

## Background

Many groups, including our own, have proposed the use of DNA methylation profiles as biomarkers for various disease states. While much research has been done identifying DNA methylation signatures in cancer vs. normal etc., we still lack sufficient knowledge of the role that differential methylation plays during normal cellular differentiation and tissue specification. We also need thorough, genome level studies to determine the meaning of methylation of individual CpG dinucleotides in terms of gene expression.

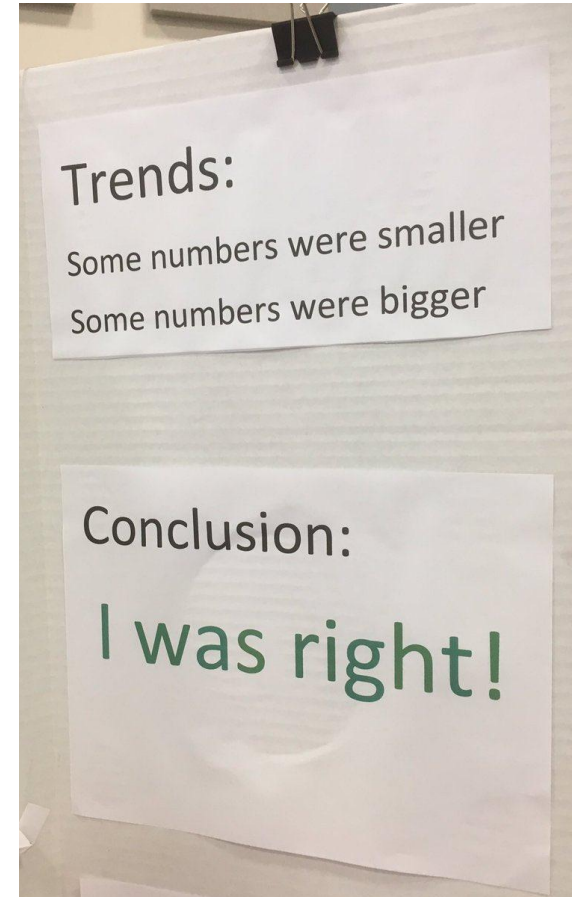
## Results

In this study, we have used (insert statistical method here) to compile unique DNA methylation signatures from normal human heart, lung, and kidney using the Illumina Infinium 27 K methylation arrays and compared those to gene expression by RNA sequencing. We have identified unique signatures of global DNA methylation for human heart, kidney and liver, and showed that DNA methylation data can be used to correctly classify various tissues. It indicates that DNA methylation reflects tissue specificity and may play an important role in tissue differentiation. The integrative analysis of methylation and RNA-Seq data showed that gene methylation and its transcriptional levels were comprehensively correlated. The location of methylation markers in terms of distance to transcription start site and CpG island showed no effects on the regulation of gene expression by DNA methylation in normal tissues.

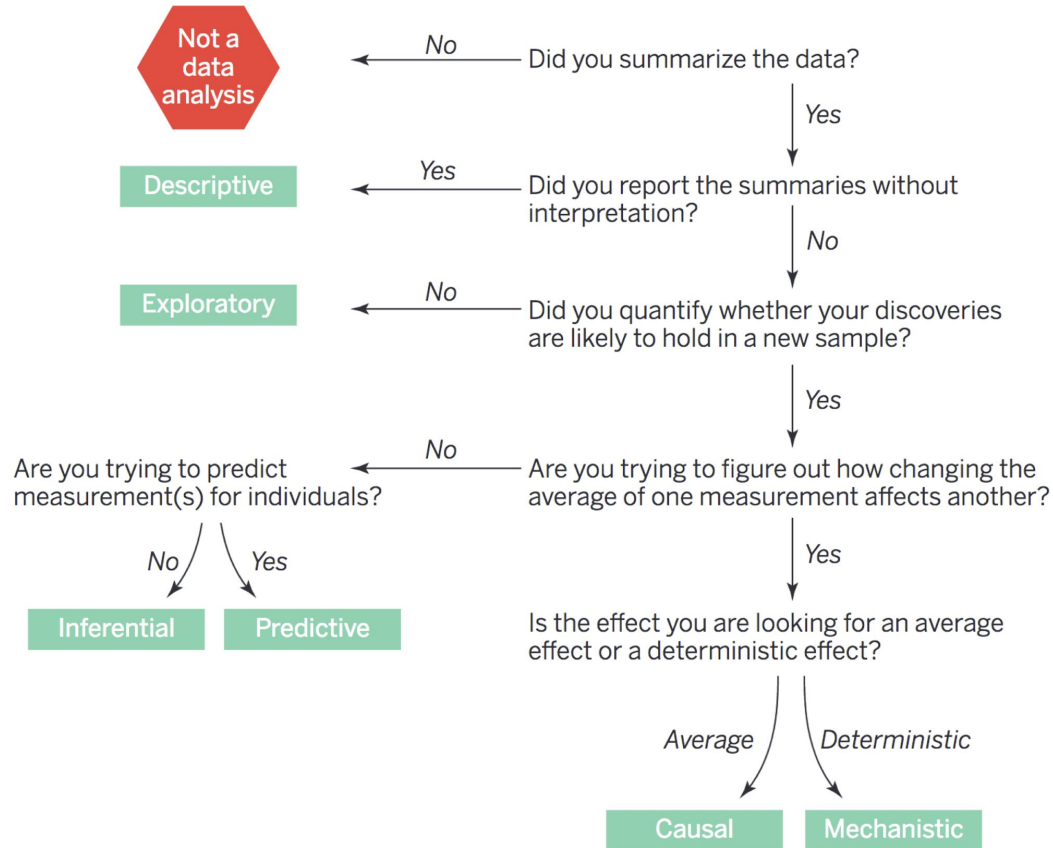
## Conclusions

This study showed that an integrative analysis of methylation array and RNA-Seq data can be utilized to discover the global regulation of gene expression by DNA methylation and suggests that DNA methylation plays an important role in normal tissue differentiation via modulation of gene expression.

<https://nsaunders.files.wordpress.com/2012/07/bmcsysbiol.png>



# Statistical hypothesis testing



# Statistical hypothesis testing

1. **Decide on the effect** that you are interested in, design a suitable experiment or study, pick a data summary function and test statistic.
2. **Set up a null hypothesis**, which is a simple, computationally tractable model of reality that lets you compute the null hypothesis.
3. **Decide on the rejection region**, i.e., a subset of possible outcomes whose total probability is small.
4. **Do the experiment** and collect the data, compute the test statistic.
5. **Make a decision**: reject the null hypothesis – i.e. conclude that it is unlikely to be true – if the test statistic is in the rejection region.

# Statistical hypothesis testing

- Many scientific studies are interested in quantifying the difference in a particular parameter between two groups.
  - There's always some difference → Is it statistically significant difference?
- **Null hypothesis:**
  - Typically a statement of no relationship between variables or no effect of an experimental manipulation/intervention.
- **Null distribution:**
  - The set of possible outcomes of the test statistic (and their probabilities) under the assumption that the null hypothesis is true.

# Statistical hypothesis testing

**Sherry:** I just love making people laugh.

I think humor is like medicine!

**Niles:** [aside] Oh, we must be in the placebo group.



- **Two groups:** Yes / No #LOL
- **Measurement:** Difference in “healthy” days
- **Null:** Ineffective | **Alternative:** Effective

# Statistical hypothesis testing

**Sherry:** I just love making people laugh.

I think humor is like medicine!

**Niles:** [aside] Oh, we must be in the placebo group.



Did you regularly #LOL more this past year than the year before?

☐ Yes

☐ No

# Statistical hypothesis testing

	Raw Data		Means
Did you regularly #LOL more this past year than the year before?	No (group1)	0, 5, 2, -2, 8, -6, 0, 0, -6, -3, -7, 4, -2, -2, 0, -4, -1, -8, 0, 6	-0.80
	Yes (group2)	-2, 2, -2, -5, 7, 11, 6, 2, 1, 1, 6, -2, -1, 7, 7, -2, 4, 3, -1, 17	2.95



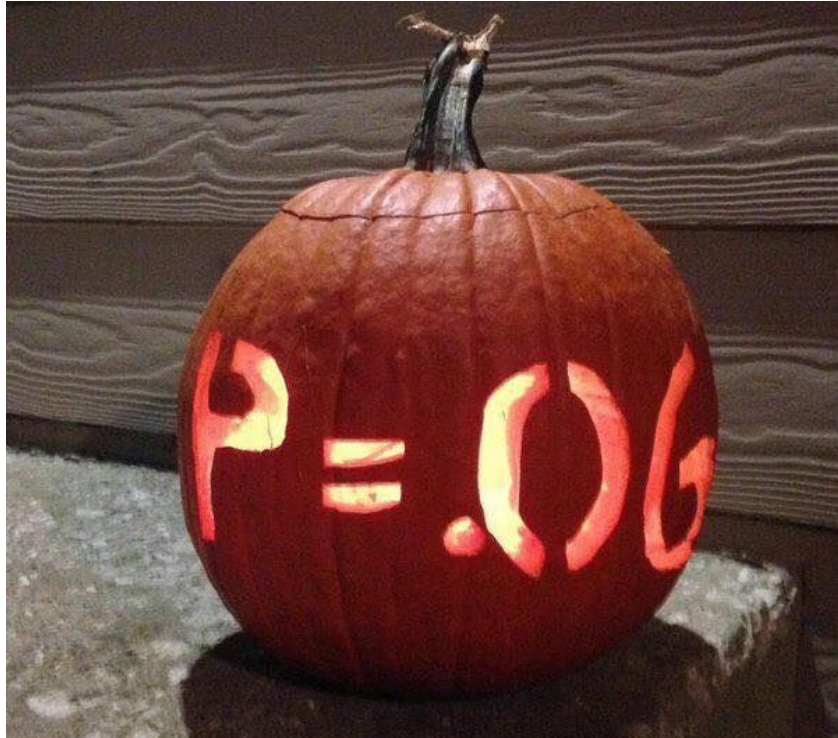
Test the efficacy  
of humor as  
medicine

Quantifying the difference between the groups	
Effect size   Diff. b/w the means of two groups ( $\mu_2 - \mu_1$ )	3.75
Test statistic   ( $\mu_2 - \mu_1$ ) / $\sqrt{(s_1^2/n_1 + s_2^2/n_2)}$	2.45



# Typical next step: perform a statistical test and get a P-value

The 'p' in p-value *actually* stands for p-potentially interesting!



ALTBIER - 4.9% ABV

The original amber ale as created by the Germans. Slightly drier than the American version, this beer drinks easy and satisfies the palette with notes of toffee and caramel without being thick or too dark which makes it a good idea.

## **P-VALUE**

DRY-HOPPED AMERICAN PALE ALE - 5.4% ABV

This Pale Ale is light and hoppy with just the right amount of malt depth. This beer challenges the notion that hops and grain can't be balanced. Reject the null hypothesis.

## **SENSORY OVERLOAD**

NEW ENGLAND IPA - 6.1% ABV

Sensory Overload doesn't let bitterness get in the way as your senses go into overdrive trying to keep up with the juicy citrus and fruit flavors we've extracted from the hops.

# What is the P-value?

- The p-value is:
  - The amount of evidence that there is an effect
  - The strength of the effect
  - The probability that the observed outcome is important
  - The probability that the intervention is ineffective

The p-value is the probability that the experiment would have produced the observed outcome – or something more extreme – if the intervention were ineffective.

# Statistical hypothesis testing

	Raw Data		Means
Did you regularly #LOL more this past year than the year before?	No (group1)	0, 5, 2, -2, 8, -6, 0, 0, -6, -3, -7, 4, -2, -2, 0, -4, -1, -8, 0, 6	-0.80
	Yes (group2)	-2, 2, -2, -5, 7, 11, 6, 2, 1, 1, 6, -2, -1, 7, 7, -2, 4, 3, -1, 17	2.95

**Null hypothesis:**  
The intervention is ineffective.

- Even if there's no intervention, we'll observe as much or more difference between the groups.
- Irrespective of which individual is part of which group, we'll observe as much or more difference b/w the groups.

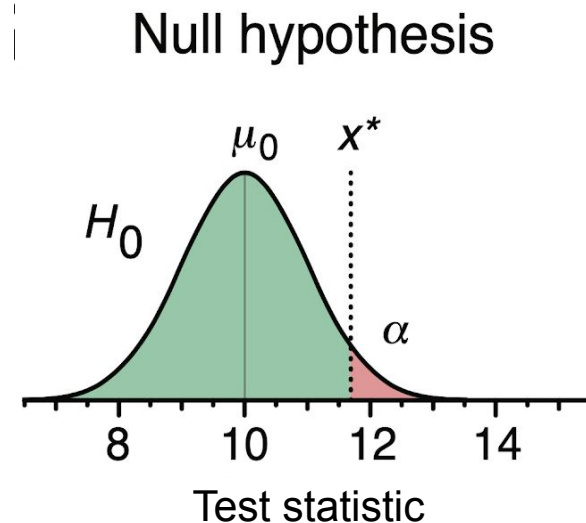


Quantifying the difference between the groups	
Effect size   Diff. b/w the means of two groups ( $\mu_2 - \mu_1$ )	1.35
Test statistic   ( $\mu_2 - \mu_1$ ) / $\sqrt{(s_1^2/n_1 + s_2^2/n_2)}$	2.45

# P-value

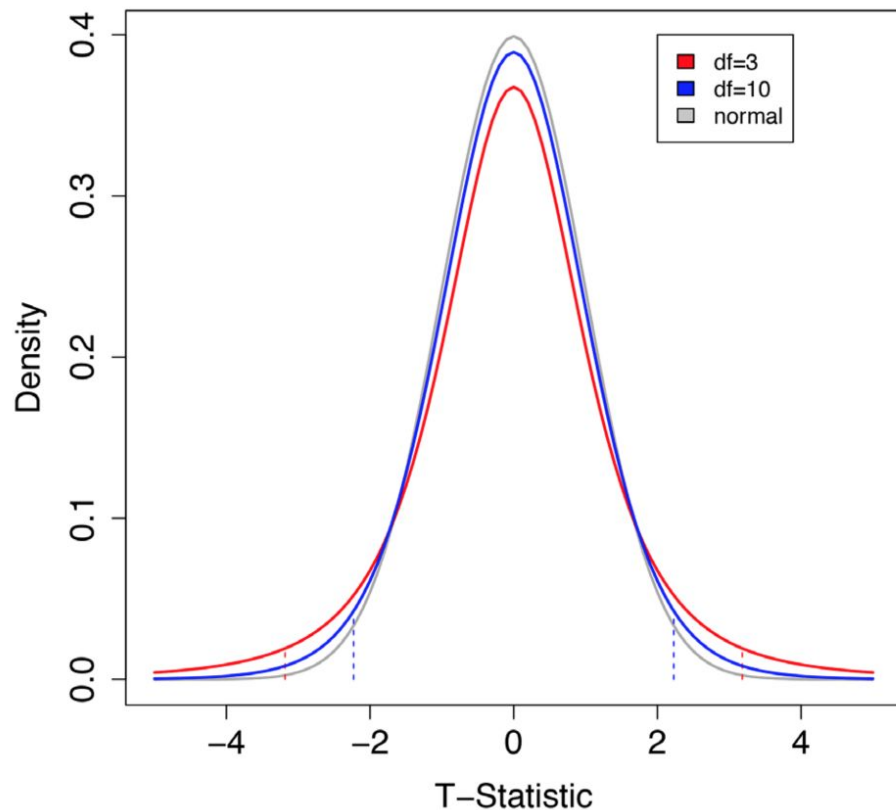
The p-value is the probability that the experiment would have produced the observed outcome – or something more extreme – if the intervention were ineffective.

1. Calculate the real test statistic.
2. Repeat the following 10,000 times to set up the null hypothesis for this test statistic:
  - Randomly assign individuals to groups.
  - Record the test statistic of the permuted data.
3. Calculate the p-value. [How?]



# P-value

T Distribution



The p-value is the area under the null distribution corresponding to outcome equal to or more extreme than the observed statistic.

Student's  
one-sample  
test

$$t = \frac{\bar{X} - \mu_0}{\text{SEM}}$$

Welch's  
two-sample  
test

$$t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{s_1^2}{N_1} + \frac{s_2^2}{N_2}}}$$

# P-value - History

- Fisher (1920s):
  - Informal method to help interpret the data along with prior experience, domain knowledge, size of the effect, etc.
- Neyman & Pearson:
  - Control false positive rate at  $\alpha$ , set by the experimenter based on what can be tolerated.
  - Formulate null and alternative hypothesis.
  - Reject null when  $p < \alpha$ .
    - The threshold  $\alpha = 0.05$  is merely a convention.

# Type I & type II errors

P-value captures if there is “sufficient” inconsistency with the null hypothesis.

Choosing  $p < \alpha$  controls type I error at  $\alpha$ .

- Type I error: False-positive rate ( $\alpha$ )
- Type II error: False-negative rate ( $\beta$ )
- Remember the story of the boy that cried wolf!



David Robinson

@drob

Follow

Remember, mixing up Type I and Type II errors is called a Type III error



David Robinson

@drob

Follow

Giving mistakes numbers instead of names was a real Type IV error

# P-value

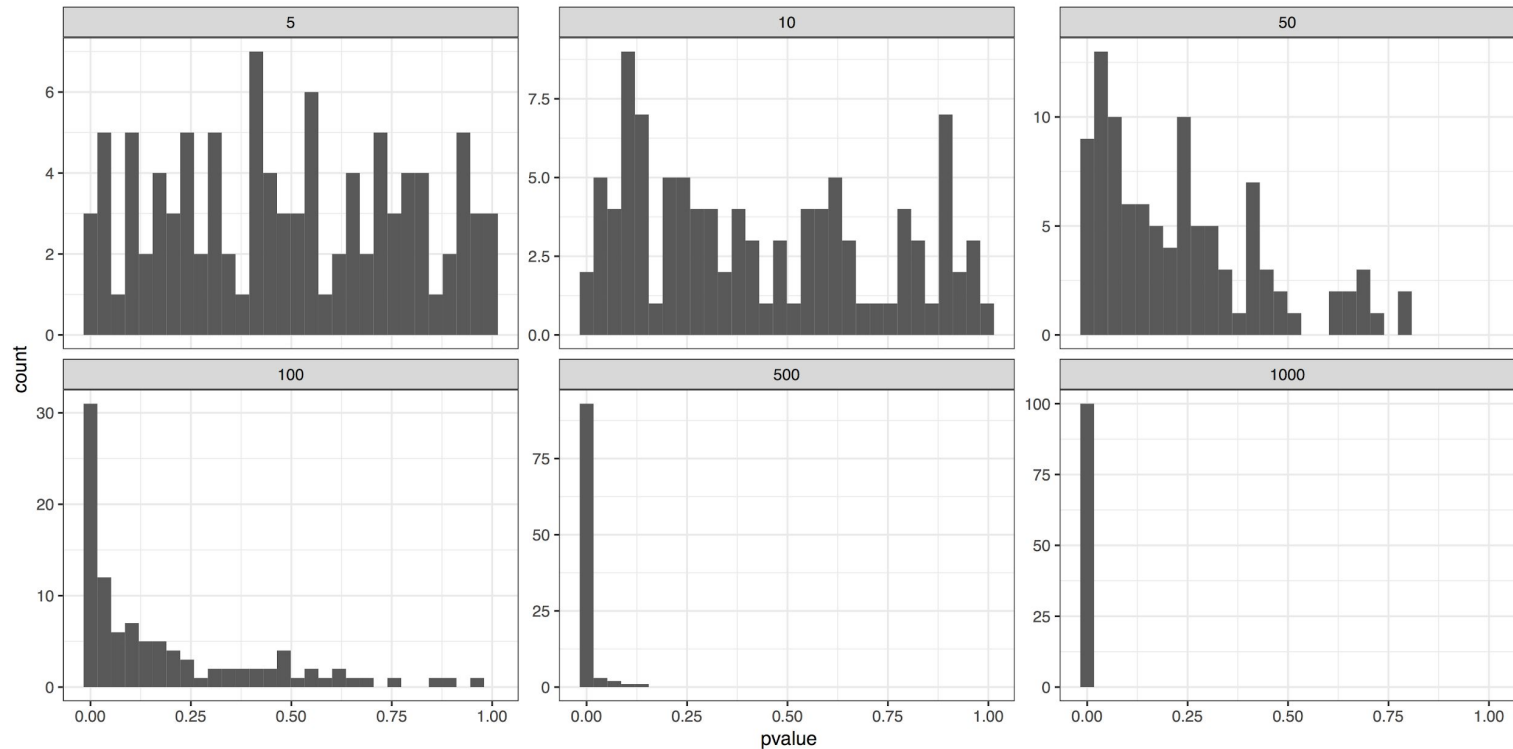
P-values are dependent on:

- Size of the effect (effect size)
- Variance within each group
- Sample size
- The underlying experimental design & the null hypothesis (need not always be random chance).
  - a. Conversely, two completely different experiments can give same data but end up very different p-values.
    - 3 out of 9: Binomial p-value = 0.073
    - 3 out of 9: Neg. Binomial p-value = 0.033.



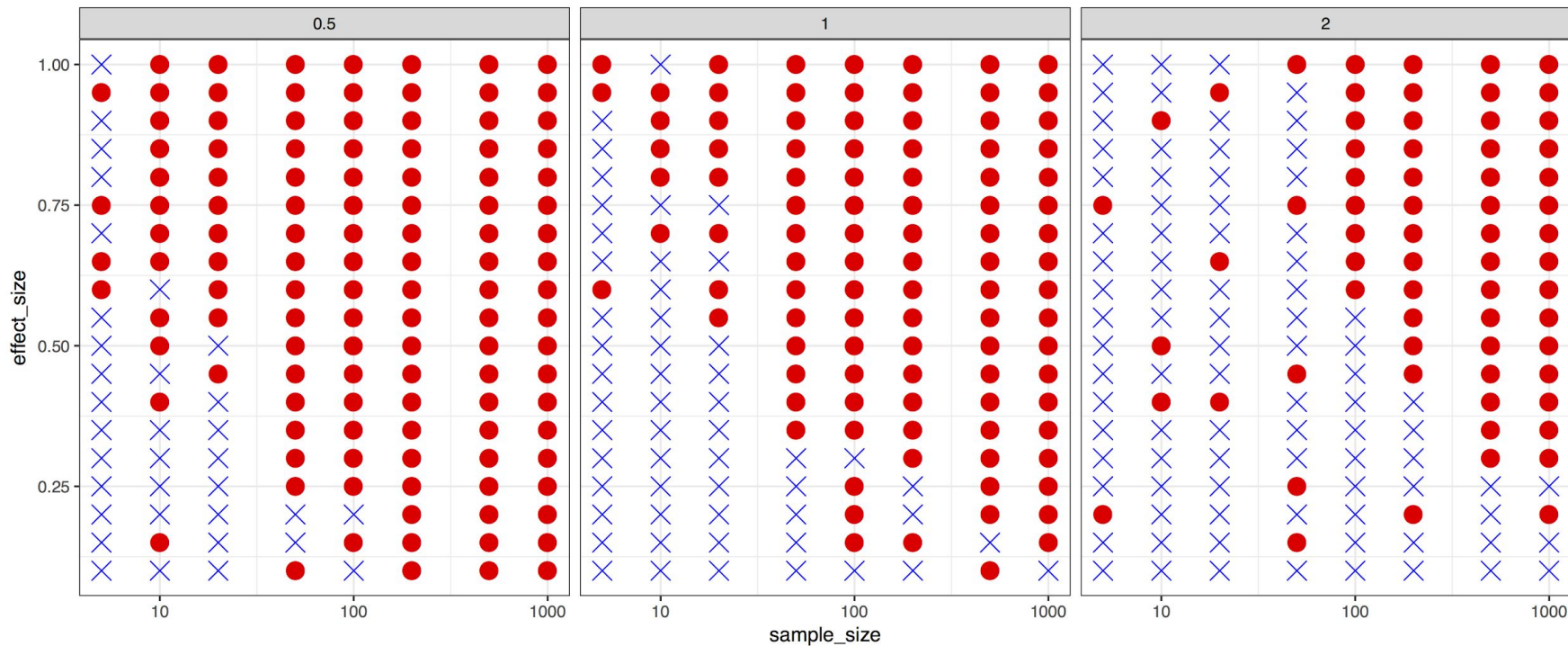
# P-value

- P-values are dependent on: sample\_size (effect\_size = 0.25, std\_deviation = 1)



# P-value

- P-values are dependent on: sample\_size, effect\_size, within-group variance



## Significant or not!

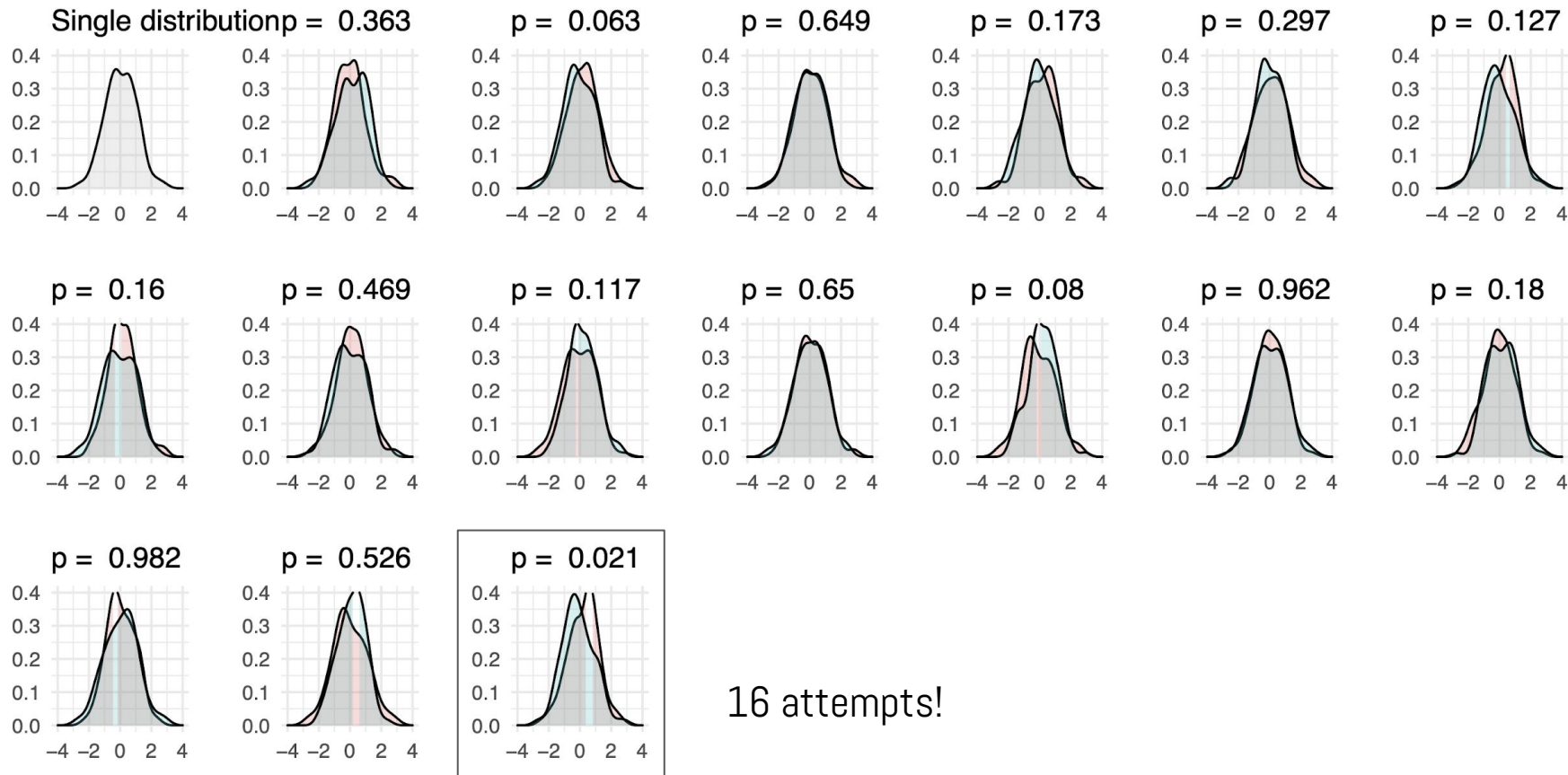
<https://mchankins.wordpress.com/2013/04/21/still-not-significant-2/>

The following list is culled from peer-reviewed journal articles in which:

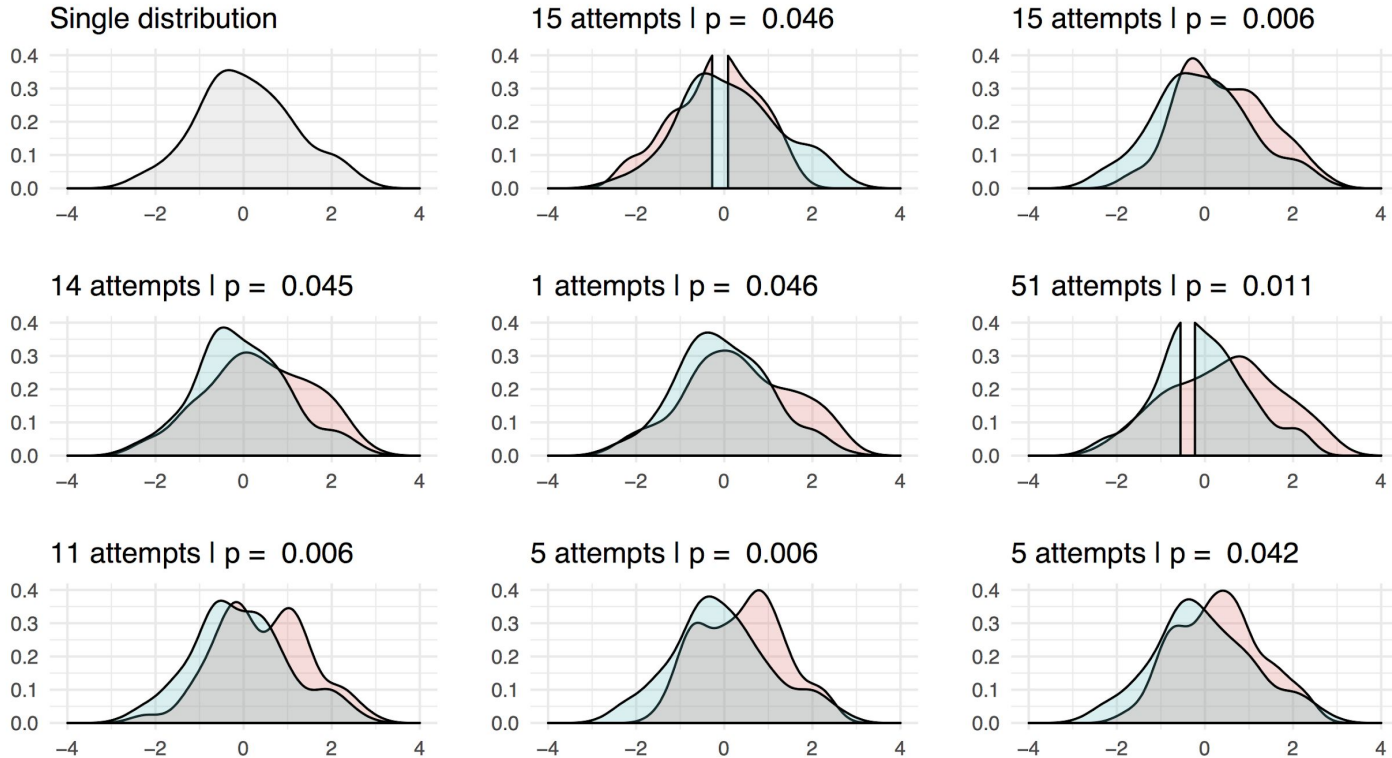
- (a) the authors set themselves the threshold of 0.05 for significance,
- (b) failed to achieve that threshold value for  $p$  and
- (c) described it in such a way as to make it seem more interesting.

(barely) not statistically significant ( $p=0.052$ )  
a barely detectable statistically significant difference ( $p=0.073$ )  
a borderline significant trend ( $p=0.09$ )  
a certain trend toward significance ( $p=0.08$ )  
a clear tendency to significance ( $p=0.052$ )  
a clear trend ( $p<0.09$ )  
a clear, strong trend ( $p=0.09$ )  
a considerable trend toward significance ( $p=0.069$ )  
a decreasing trend ( $p=0.09$ )  
a definite trend ( $p=0.08$ )  
a distinct trend toward significance ( $p=0.07$ )  
a favorable trend ( $p=0.09$ )

# P-hacking and Publication bias

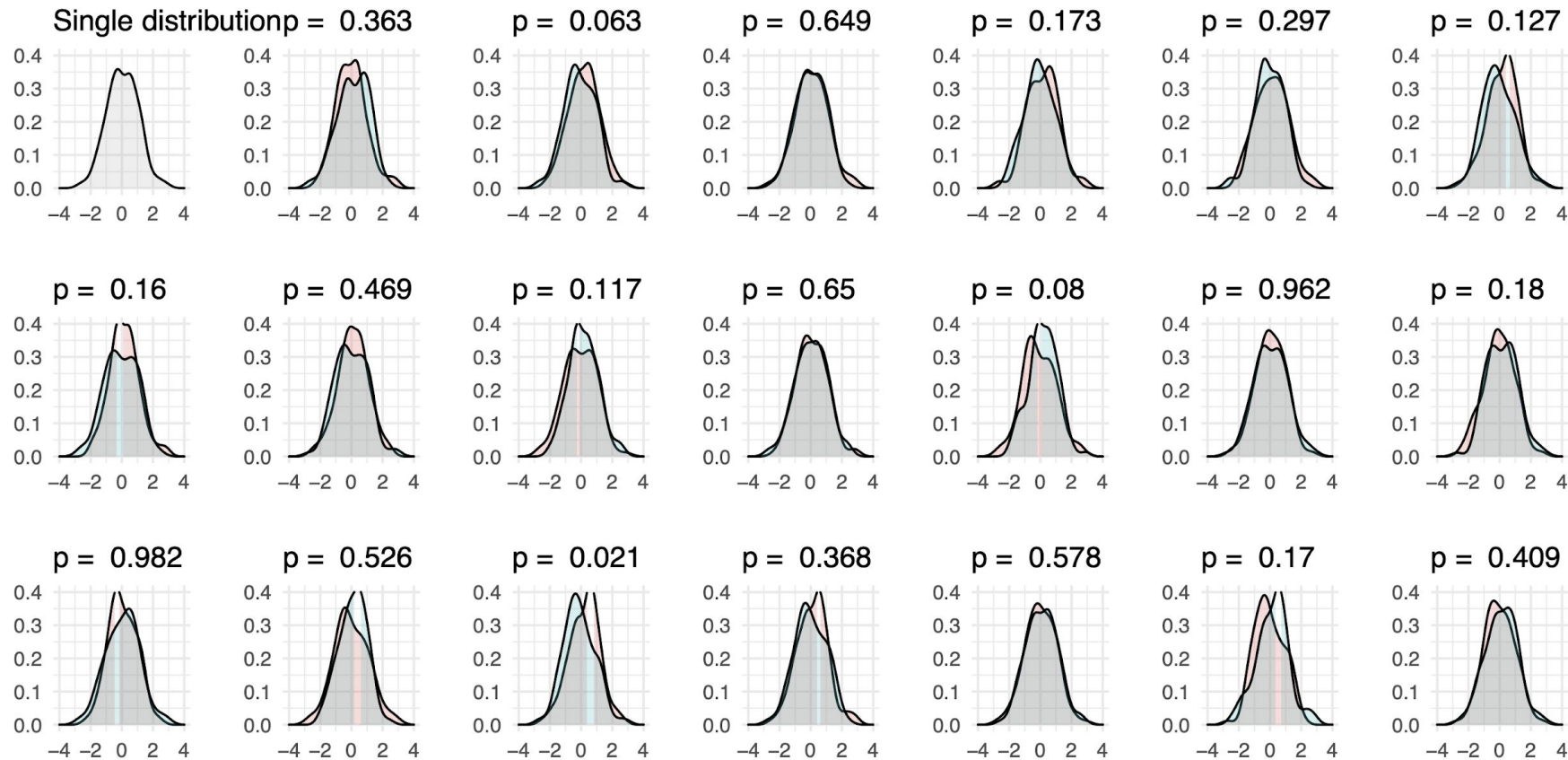


# P-hacking and Publication bias



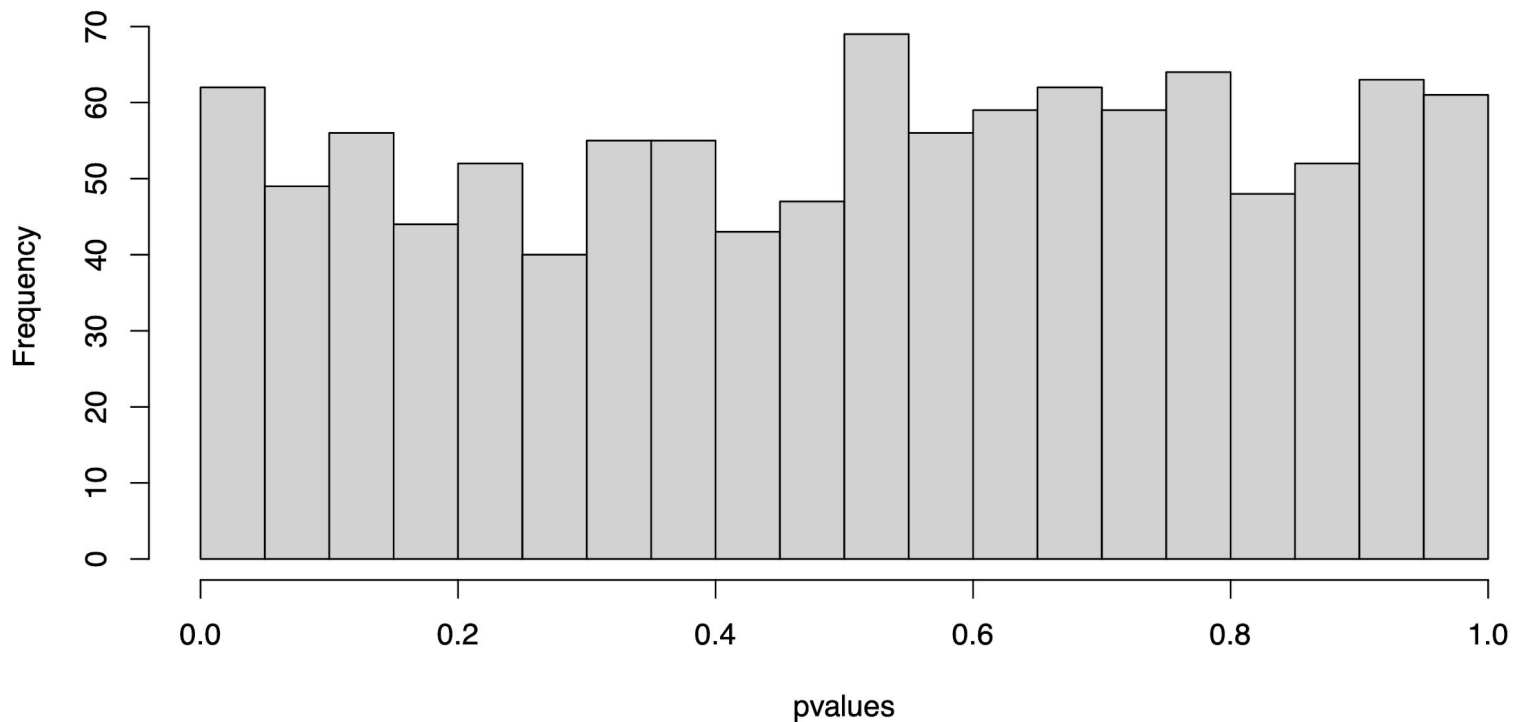
“When a measure become a target, it ceases to be a good measure” – Goodhart's Law

# P-hacking and Publication bias



# P-hacking and Publication bias

Distribution of p-values under the null hypothesis



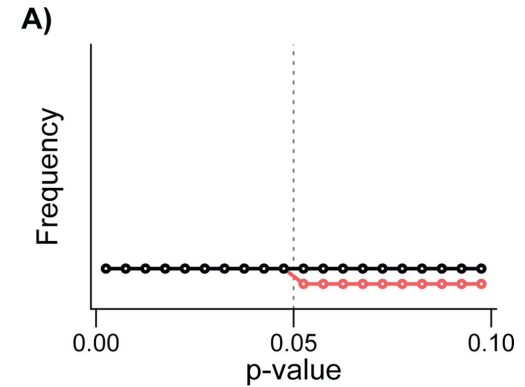
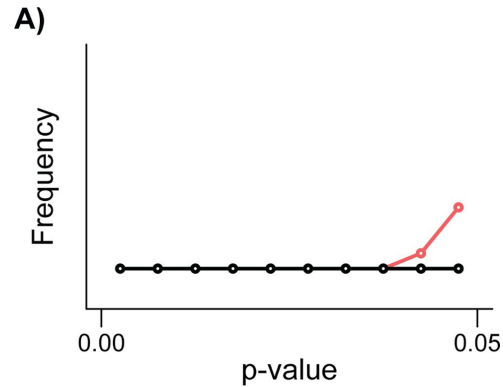
# P-hacking and Publication bias

- **P-hacking:** Collect or select data or statistical analyses until nonsignificant results become significant.
  - Conducting analyses midway through experiments to decide whether to continue collecting data.
  - Recording many response variables and deciding which to report post-analysis
  - Deciding whether to include or drop outliers post-analyses
  - Excluding, combining, or splitting treatment groups post-analysis
  - Including or excluding covariates post-analysis, and
  - Stopping data exploration if an analysis yields a significant p-value.
- **Publication bias:** Studies with nonsignificant results have lower publication rates.

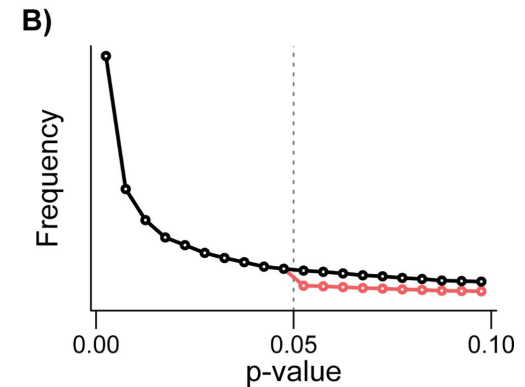
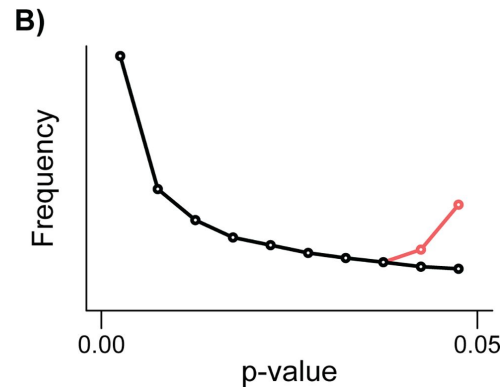


# P-hacking and Publication bias

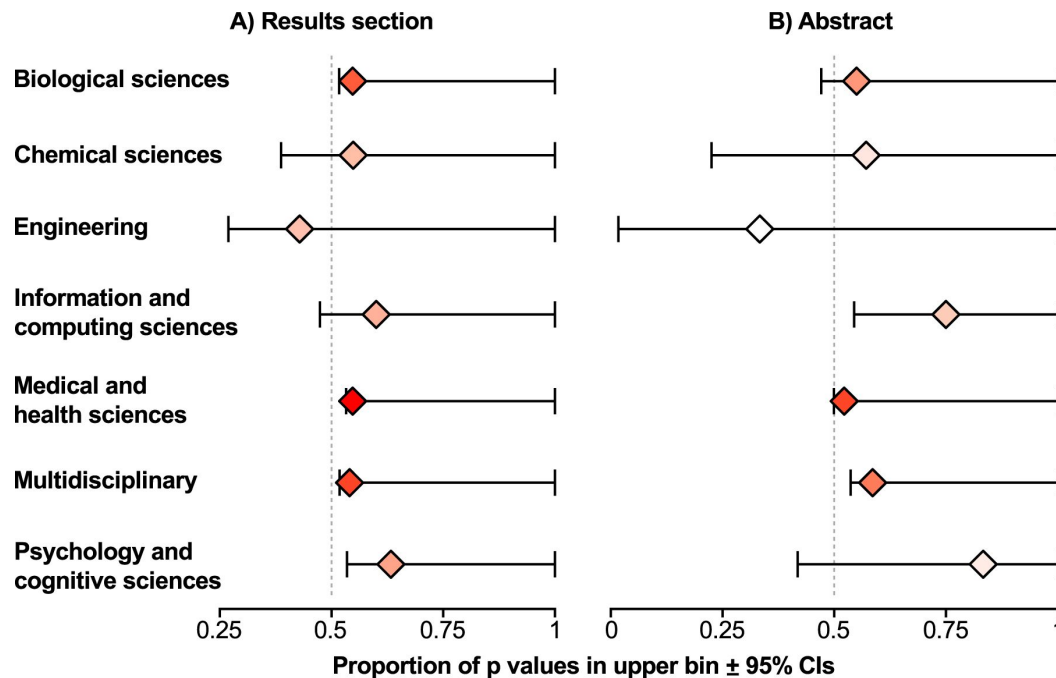
No evidence  
of effect



Evidence  
of effect



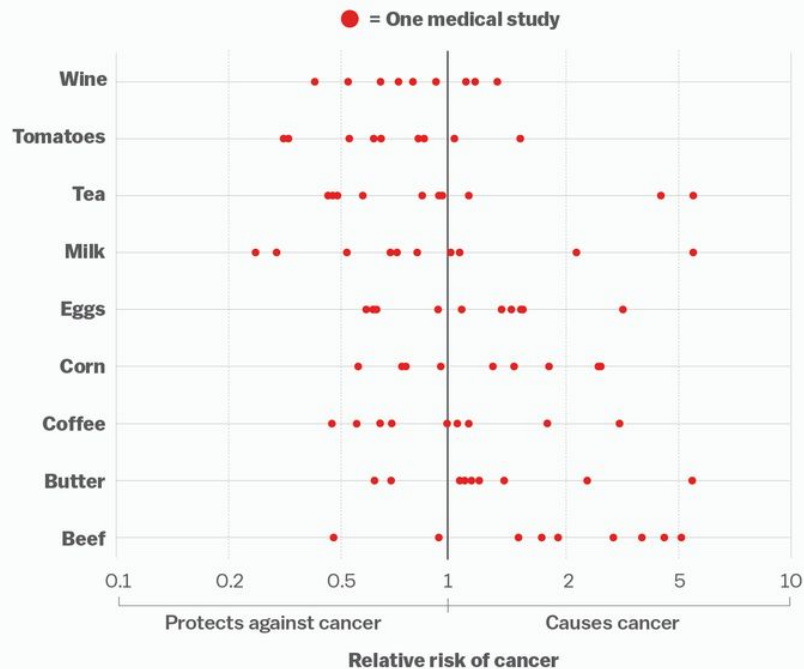
# Publication bias and P-hacking



Hack Your Way To Scientific Glory: <https://projects.fivethirtyeight.com/p-hacking/>

# Statistical hypothesis testing

Everything we eat both causes and prevents cancer



SOURCE: Schoenfeld and Ioannidis, *American Journal of Clinical Nutrition*

Vox



TIME  
@TIME

How coffee can help you live longer



How Coffee Can Help You Live Longer  
New findings add to growing evidence that co...  
time.com

4/9/17, 6:45 AM



TIME  
@TIME

The problem with your coffee



Hot Drinks a Probable Cancer Cause, Says WHO  
time.com

4/9/17, 6:15 AM

# Questionable research practices

- Exclusively using p-values to determine the relevance and sanity of the results of a statistical test.
- Analyzing the data until the desired results are found.
- Collecting more data to reach smaller p-values.
- Trying many hypothesis until one of them gives a low p-value, and reporting just that final result.

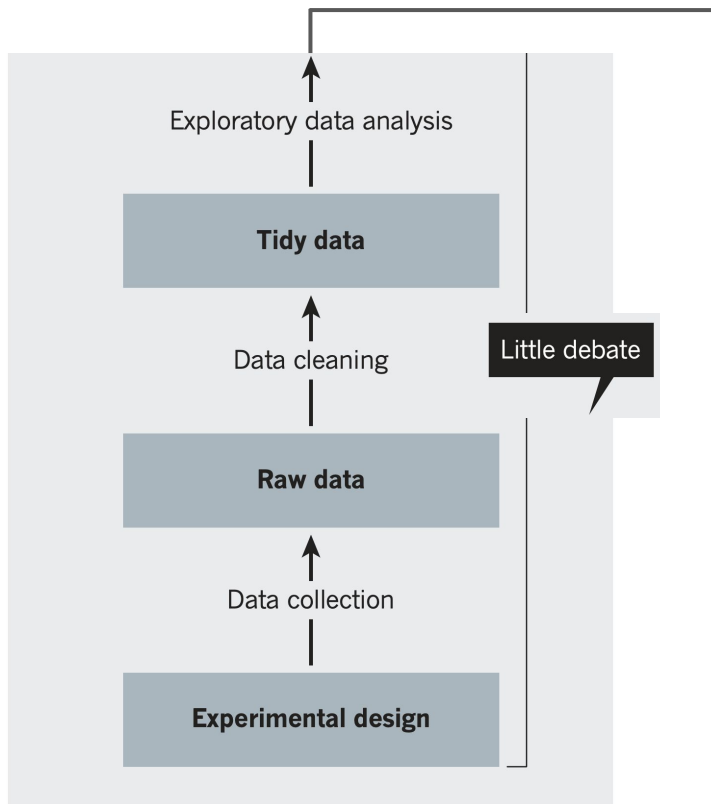
WHEN YOU SEE A CLAIM THAT A COMMON DRUG OR VITAMIN "KILLS CANCER CELLS IN A PETRI DISH,"

KEEP IN MIND:



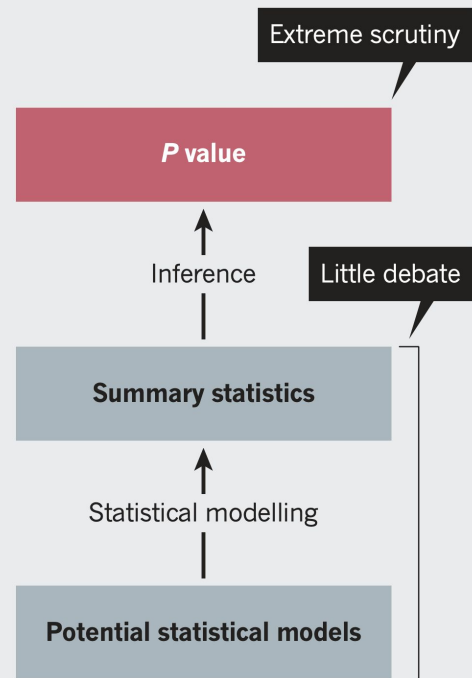
SO DOES A HANDGUN.

# P-values are just the tip of the iceberg!

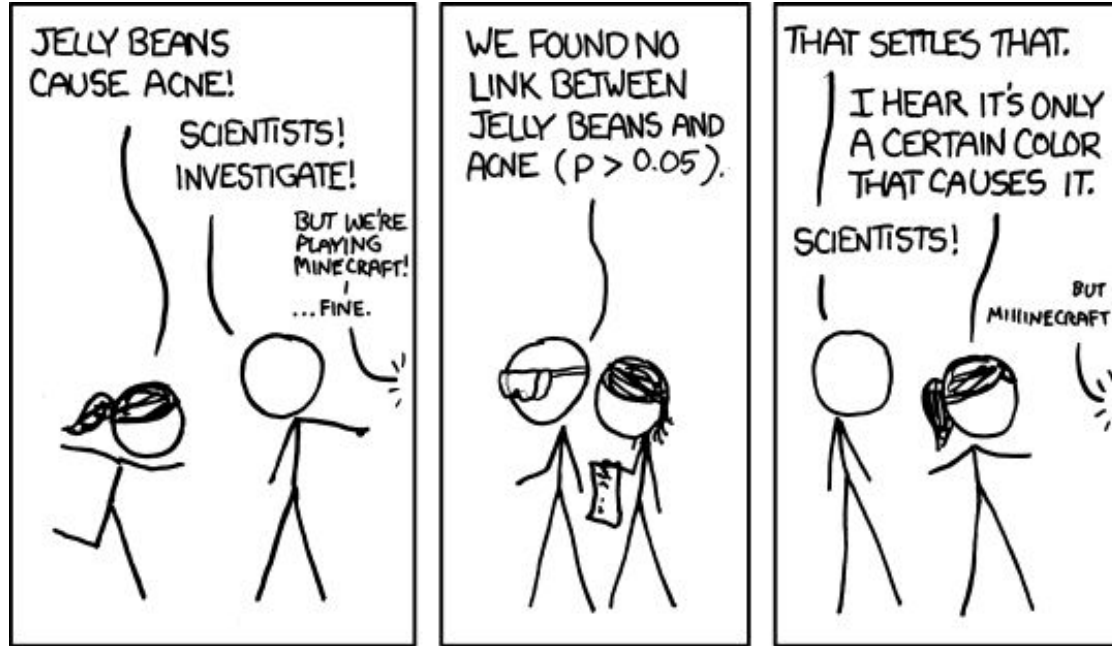


## DATA PIPELINE

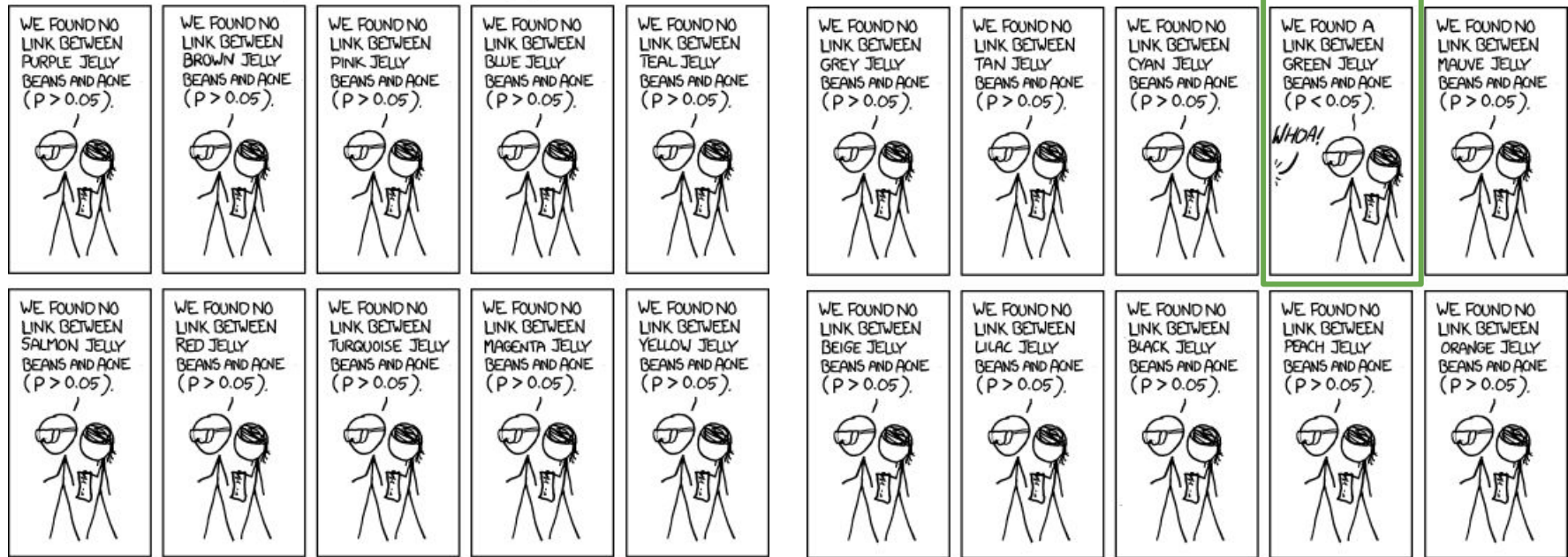
The design and analysis of a successful study has many stages, all of which need policing.



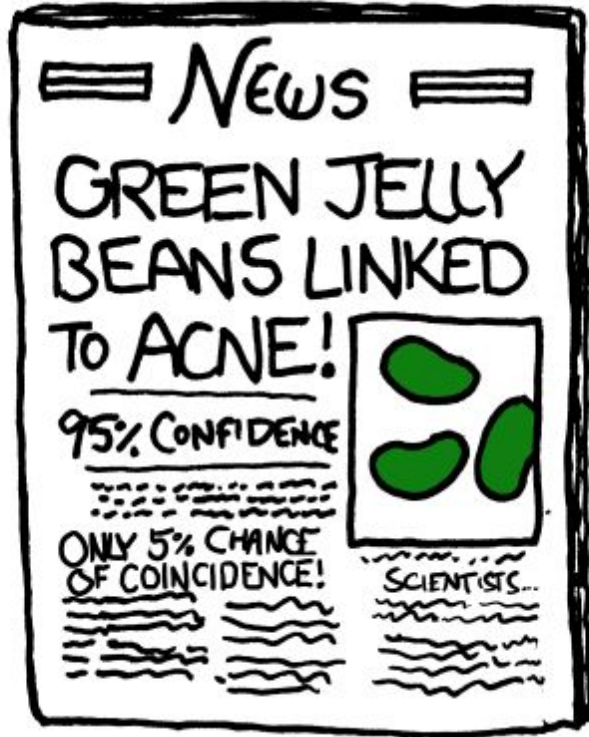
# Multiple hypothesis testing



# Multiple hypothesis testing



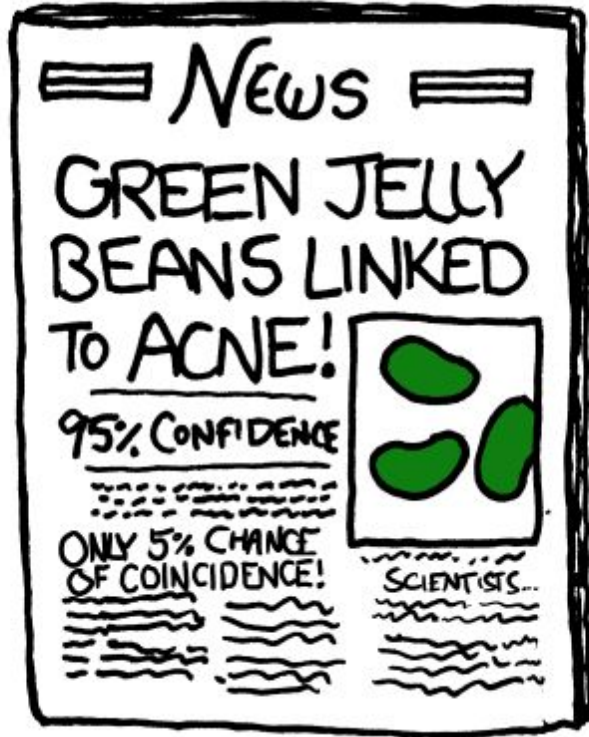
# Multiple hypothesis testing



- The more inferences are made, the more likely erroneous inferences are to occur.
- This issue of multiple testing is related to:
  - P-hacking
  - Publication bias



# Multiple hypothesis testing



Several statistical techniques have been developed to prevent this from happening.

- These techniques generally require a **stricter significance threshold for individual comparisons**, so as to compensate for the number of inferences being made.
- We are going to discuss later the idea of **False Discovery Rate**.

# What you need to do before the next class

- Complete the assignment
  - Implementing the permutation test
  - Exploring the dependence of p-value on effect size, sample size, & variance
- Concepts
  - Brush-up: Statistical power
  - Brush-up: Replication