# Statistical Significance in Research

Thomas.Love@case.edu

2020-11-16

# Where To Get More Information

- Today's slides are at https://github.com/THOMASELOVE/rethink.
- You'll find all of the references there, as well, and some other sources.

**Thomas E. Love, Ph.D.**

- Professor of Medicine, Population & Quantitative Health Sciences, CWRU School of Medicine
- Director of Biostatistics and Data Science, Population Health Research Institute, The MetroHealth System
- Chief Data Scientist, Better Health Partnership
- Fellow, American Statistical Association

My email is Thomas dot Love at case dot edu.

# Today's Agenda

- What I Taught for Many Years
- Notions of Statistical Significance
  - *p* values
  - How the academics have reacted
- Evaluating Health News: A primer
- What I Think I Think Now

"Not Even Scientists Can Easily Explain p Values" Video
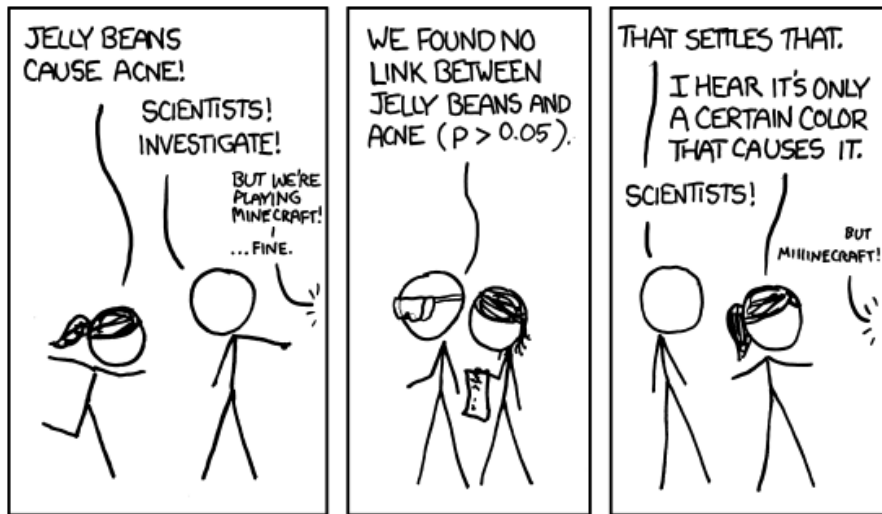
Section 1

## What I Taught for Many Years

# What I Taught for Many Years

- Null hypothesis significance testing is here to stay.
  - Learn how to present your p value so it looks like what everyone else does
  - Think about "statistically detectable" rather than "statistically significant"
  - Don't accept a null hypothesis, just retain it.
- Use point **and** interval estimates
  - Try to get your statements about confidence intervals right (right = just like I said it)
- Use Bayesian approaches/simulation/hierarchical models when they seem appropriate or for "non-standard" designs
  - But look elsewhere for people to teach/do that stuff
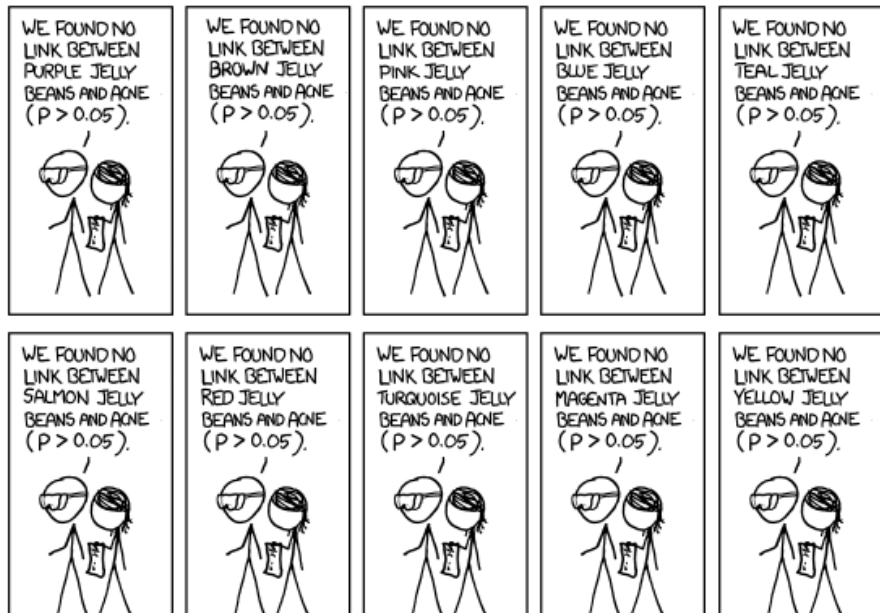- Power is basically a hurdle to overcome in a grant application

# Conventions for Reporting *p* Values

1. Use an italicized, lower-case *p* to specify the *p* value. Don't use *p* for anything else.
2. For *p* values above 0.10, round to two decimal places, at most.
3. For *p* values near $\alpha$, include only enough decimal places to clarify the reject/retain decision.
4. For very small *p* values, always report either $p < 0.0001$ or even just $p < 0.001$, rather than specifying the result in scientific notation, or, worse, as $p = 0$ which is glaringly inappropriate.
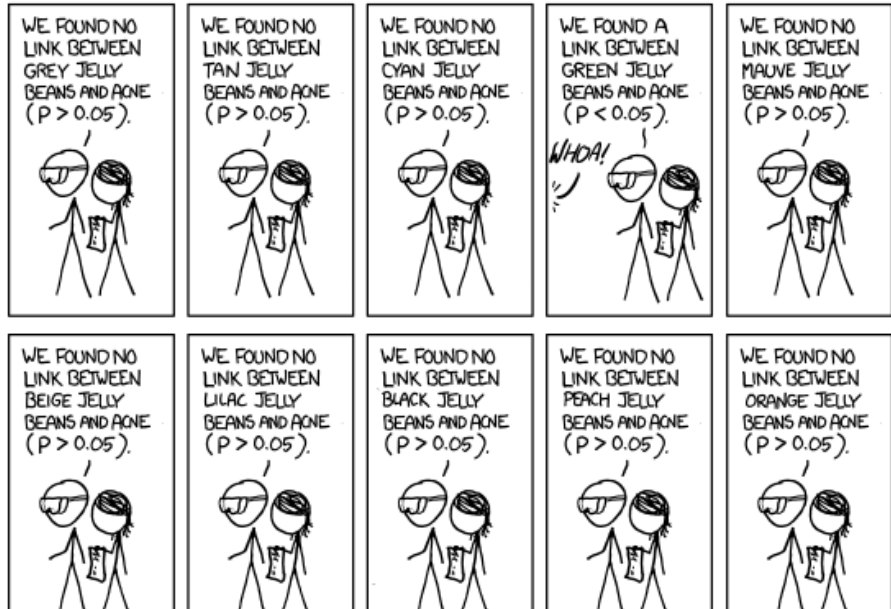5. Report *p* values above 0.99 as $p > 0.99$, rather than $p = 1$.

# XKCD "Significance"

Section 2

# American Statistical Association to the rescue!?!

# ASA 2016 Statement on *p* Values

ASA Statement: "Informally, a p-value is the probability under a specified statistical model that a statistical summary of the data (e.g., the sample mean difference between two compared groups) would be equal to or more extreme than its observed value."

*fivethirtyeight*.com "Not Even Scientists Can Easily Explain *p* Values"

. . . Try to distill the p-value down to an intuitive concept and it loses all its nuances and complexity, said science journalist Regina Nuzzo, a statistics professor at Gallaudet University. "Then people get it wrong, and this is why statisticians are upset and scientists are confused." **You can get it right, or you can make it intuitive, but it's all but impossible to do both.**

*fivethirtyeight*.com "Statisticians found one thing they can agree on"

# A Few Comments on Significance

- **A significant effect is not necessarily the same thing as an interesting effect.** For example, results calculated from large samples are nearly always "significant" even when the effects are quite small in magnitude. Before doing a test, always ask if the effect is large enough to be of any practical interest. If not, why do the test?

- **A non-significant effect is not necessarily the same thing as no difference.** A large effect of real practical interest may still produce a non-significant result simply because the sample is too small.

- **There are assumptions behind all statistical inferences.** Checking assumptions is crucial to validating the inference made by any test or confidence interval.

- "**Scientific conclusions and business or policy decisions should not be based only on whether a p-value passes a specific threshold.**"

ASA *statement* on *p* values

# From George Cobb - on why *p* values deserve to be re-evaluated

The **idea** of a p-value as one possible summary of evidence

morphed into a

- **rule** for authors: reject the null hypothesis if p < .05.

# From George Cobb - on why *p* values deserve to be re-evaluated

The **idea** of a p-value as one possible summary of evidence morphed into a

- **rule** for authors: reject the null hypothesis if $p < .05$,

which morphed into a

- **rule** for editors: reject the submitted article if $p > .05$.

# From George Cobb - on why *p* values deserve to be re-evaluated

The **idea** of a p-value as one possible summary of evidence morphed into a

- **rule** for authors: reject the null hypothesis if $p < .05$,

which morphed into a

- **rule** for editors: reject the submitted article if $p > .05$,

which morphed into a

- **rule** for journals: reject all articles that report p-values[1]

---

[1]http://www.nature.com/news/psychology-journal-bans-p-values-1.17001 describes the recent banning of null hypothesis significance testing by *Basic and Applied Psychology*.

# From George Cobb - on why *p* values deserve to be re-evaluated

The **idea** of a p-value as one possible summary of evidence

morphed into a

- **rule** for authors: reject the null hypothesis if $p < .05$, which morphed into a
- **rule** for editors: reject the submitted article if $p > .05$, which morphed into a
- **rule** for journals: reject all articles that report p-values.

Bottom line: **Reject rules. Ideas matter.**

# Roger Peng's description of a successful data analysis

*A data analysis is successful if the audience to which it is presented accepts the results.*

- "What is a Successful Data Analysis?" https://simplystatistics.org/2018/04/17/what-is-a-successful-data-analysis/.

So what makes a data analysis more believable / more acceptable?

# The American Statistical Association

2016

- Ronald L. Wasserstein & Nicole A. Lazar (2016) The ASA's Statement on p-Values: Context, Process, and Purpose, *The American Statistician*, 70:2, 129-133, DOI: 10.1080/00031305.2016.1154108

2019

- Ronald L. Wasserstein, Allen L. Schirm & Nicole A. Lazar (2019) Moving to a World Beyond "$p < 0.05$", *The American Statistician*, 73:sup1, 1-19, DOI: 10.1080/00031305.2019.1583913.

# Statistical Inference in the 21st Century

*. . . a world learning to venture beyond "p < 0.05"*

*This is a world where researchers are free to treat "p = 0.051" and "p = 0.049" as not being categorically different, where authors no longer find themselves constrained to selectively publish their results based on a single magic number.*

# Statistical Inference in the 21st Century

*In this world, where studies with "$p < 0.05$" and studies with "$p > 0.05$" are not automatically in conflict, researchers will see their results more easily replicated – and, even when not, they will better understand why.*

*The 2016 ASA Statement on P-Values and Statistical Significance started moving us toward this world. As of the date of publication of this special issue, the statement has been viewed over 294,000 times and cited over 1700 times-an average of about 11 citations per week since its release. Now we must go further.*

# The American Statistical Association Statement on P values and Statistical Significance

The ASA Statement (2016) was mostly about what **not** to do.

The 2019 effort represents an attempt to explain what to do.

# ASA 2019 Statement

*Some of you exploring this special issue of The American Statistician might be wondering if it's a scolding from pedantic statisticians lecturing you about what not to dowith p-values, without offering any real ideas of what to do about the very hard problem of separating signal from noise in data and making decisions under uncertainty. Fear not. In this issue, thanks to 43 innovative and thought-provoking papers from forward-looking statisticians, help is on the way.*

# "Don't" is not enough.

*If you're just arriving to the debate, here's a sampling of what not to do.*

- Don't base your conclusions solely on whether an association or effect was found to be "statistically significant" (i.e., the $p$ value passed some arbitrary threshold such as $p < 0.05$).
- Don't believe that an association or effect exists just because it was statistically significant.

# "Don't" is not enough.

- Don't believe that an association or effect is absent just because it was not statistically significant.
- Don't believe that your p-value gives the probability that chance alone produced the observed association or effect or the probability that your test hypothesis is true.
- Don't conclude anything about scientific or practical importance based on statistical significance (or lack thereof).

# One More Don't. . .

The *ASA Statement on P-Values and Statistical Significance* stopped just short of recommending that declarations of "statistical significance" be abandoned. We take that step here. We conclude, based on our review of the articles in this special issue and the broader literature, that it is time to stop using the term "statistically significant" entirely. Nor should variants such as "significantly different," "$p < 0.05$," and "nonsignificant" survive, whether expressed in words, by asterisks in a table, or in some other way.

Regardless of whether it was ever useful, a declaration of "statistical significance" has today become meaningless. Made

*A label of statistical significance adds nothing to what is already conveyed by the value of p; in fact, this dichotomization of p-values makes matters worse.*

Section 3

**p = 0.05?**

# p = 0.05?

*"For decades, the conventional p-value threshold has been 0.05,"*
*says Dr. Paul Wakim, chief of the biostatistics and clinical epidemi-*
*ology service at the National Institutes of Health Clinical Center,*
*"but it is extremely important to understand that this 0.05, there's*
*nothing rigorous about it. It wasn't derived from statisticians who*
*got together, calculated the best threshold, and then found that*
*it is 0.05. No, it's Ronald Fisher, who basically said, 'Let's use*
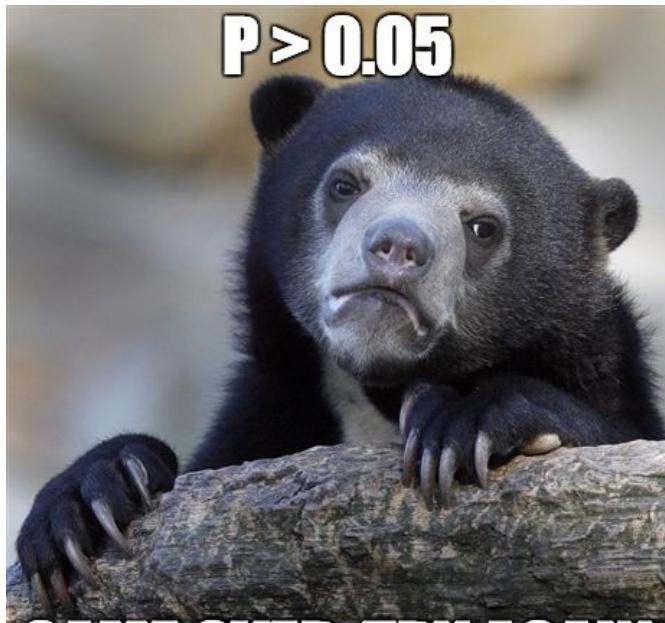*0.05,' and he admitted that it was arbitrary."*

- NOVA "Rethinking Science's Magic Number" by Tiffany Dill
  2018-02-28. See especially the video labeled "Science's most
  important (and controversial) number has its origins in a British
  experiment involving milk and tea."

# More from Dr. Wakim. . .

*"People say, 'Ugh, it's above 0.05, I wasted my time.' No, you didn't waste your time."* says Dr. Wakim. *"If the research question is important, the result is important. Whatever it is."*

- NOVA Season 45 Episode 6 Prediction by the Numbers 2018-02-28.

**Randy Sweis, MD**
@RandySweisMD

**Follow**

If a P value of 0.06 trends toward statistical significance, then doesn't a P value of 0.04 trend toward non-significance?

9:47 AM - 12 Jan 2018

# George Cobb's Questions (with Answers)

In February 2014, George Cobb, Professor Emeritus of Mathematics and Statistics at Mount Holyoke College, posed these questions to an ASA discussion forum:

Q: Why do so many colleges and grad schools teach $p = 0.05$?

A: Because that's **still** what the scientific community and journal editors use.

Q: Why do so many people still use $p = 0.05$?

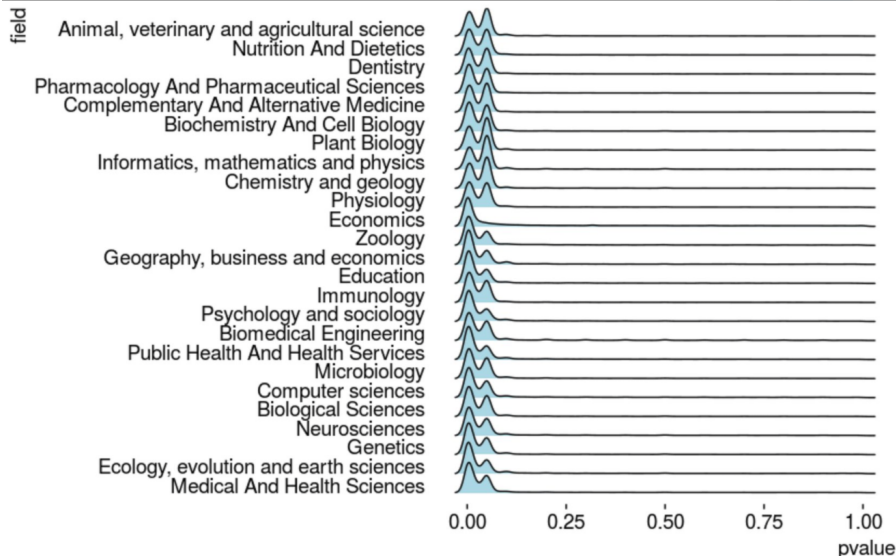A: Because that's what they were taught in college or grad school.

# All the p values

*The p-value is the most widely-known statistic. P-values are reported in a large majority of scientific publications that measure and report data. R.A. Fisher is widely credited with inventing the p-value. If he was cited every time a p-value was reported his paper would have, at the very least, 3 million citations - making it the most highly cited paper of all time.*

- Visit Jeff Leek's Github for tidypvals package
  - 2.5 million *p* values in 25 scientific fields

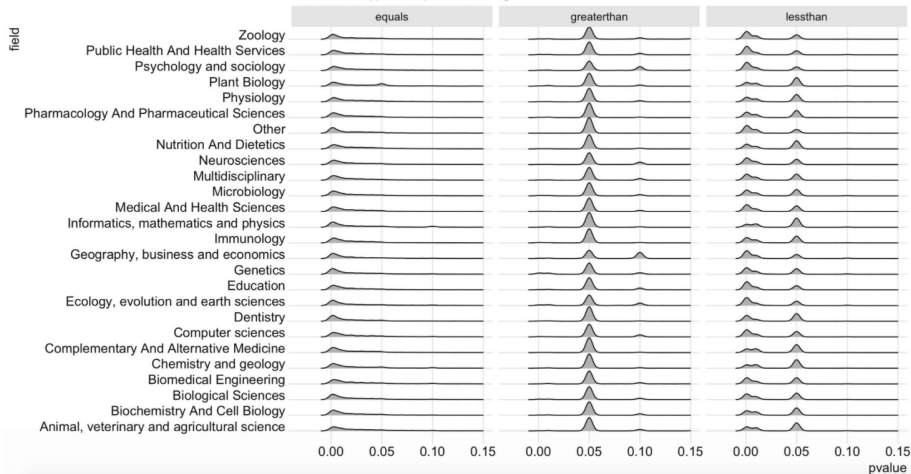**What do you suppose the distribution of those p values is going to look like?**

**Distribution of pvalues by operator (=, >, <)**
Economics dropped: all operators missing

Section 4

# The Academy Reacts

**Mustafa Ascha**                    Nov 22 (6 days ago)

to me

Hi Dr Love,

I wanted to share: JAMA-Otolaryngology just asked for a revision of a manuscript I'm writing. They specifically asked that I remove all p-values from the paper, and attached a statement from the journal editor that essentially repeated the ASA statement on p-values.

I'm pleasantly surprised by their request.

http://jamanetwork.com/journals/jamaotolaryngology/fullarticle/2546529

# Improving the Quality of the Reporting of Research Results

Jay F. Piccirillo, MD[1]

≫ Author Affiliations | Article Information

# Problems with *P* Values

1. *P* values are inherently unstable
2. The *p* value, or statistical significance, does not measure the size of an effect or the importance of a result
3. Scientific conclusions should not be based only on whether a *p* value passes a specific threshold
4. Proper inference requires full reporting and transparency
5. By itself, a *p* value does not provide a good measure of evidence regarding a model or hypothesis

http://jamanetwork.com/journals/jamaotolaryngology/fullarticle/2546529

# Solutions to the *P* Value Problems

1. Estimation of the Size of the Effect
2. Precision of the Estimate (Confidence Intervals)
3. Inference About the Target Population
4. Determination of Whether the Results Are Compatible With a Clinically Meaningful Effect
5. Replication and Steady Accumulation of Knowledge

http://jamanetwork.com/journals/jamaotolaryngology/fullarticle/2546529

# Importance of Meta-Analytic Thinking

*In JAMA Otolaryngology: Head & Neck Surgery, we look to publish original investigations where the investigators planned the study with sufficient sample size to have adequate power to detect a clinically meaningful effect and report the results with effect sizes and CIs. Authors should interpret the effect sizes in relation to previous research and use CIs to help determine whether the results are compatible with clinically meaningful effects. And finally, we acknowledge that no single study can define truth and that the advancement of medical knowledge and patient care depends on the steady accumulation of reliable clinical information.*

http://jamanetwork.com/journals/jamaotolaryngology/fullarticle/2546529

# The Value of a *p*-Valueless Paper

Jason T. Connor (2004) *American J of Gastroenterology* 99(9): 1638-40.

Abstract: As is common in current bio-medical research, about 85% of original contributions in *The American Journal of Gastroenterology* in 2004 have reported *p*-values. However, none are reported in this issue's article by Abraham et al. who, instead, rely exclusively on effect size estimates and associated confidence intervals to summarize their findings. **Authors using confidence intervals communicate much more information in a clear and efficient manner than those using *p*-values. This strategy also prevents readers from drawing erroneous conclusions caused by common misunderstandings about *p*-values**. I outline how standard, two-sided confidence intervals can be used to measure whether two treatments differ or test whether they are clinically equivalent.

*Link*

--- Editor's Note ---

## Do Not Over (*P*) Value Your Research Article

Laine E. Thomas, PhD; Michael J. Pencina, PhD

***P* value** is by far the most prevalent statistic in the medical literature but also one attracting considerable controversy. Recently, the American Statistical Association[1] released a policy statement on *P* values, noting that misunderstanding and misuse of *P* values is an important contributing factor to the common problem of scientific conclusions that fail to be reproducible. Furthermore, reliance on *P* values may distract from the good scientific principles that are needed for high-quality research. Mark et al[2] delve deeper into the history and interpretation of the *P* value in this issue of *JAMA Cardiology*. Herein, we take the opportunity to state a few principles to help guide authors in the use and reporting of *P* values in the journal.

**Related article**

When the limitations surrounding *P* values are emphasized, a common question is, "What should we do instead?" Ron Wasserstein of the American Statistical Association explained: "In the post p<0.05 era, scientific argumentation is not based on whether a p-value is small enough or not. Attention is paid to effect sizes and confidence intervals. Evidence is thought of as being continuous rather than some sort of dichotomy.... Instead, journals [should evaluate] papers based on clear and detailed description of the study design, execution, and analysis, having conclusions that are based on valid

We suggest that researchers submitting manuscripts to *JAMA Cardiology* should also consider the following:

1. Data that are descriptive of the sample (ie, indicating imbalances between observed groups but not making inference to a population) should not be associated with *P* values. Appropriate language, in this case, would describe numerical differences and sample summary statistics and focus on differences of clinical importance.

2. In addition to summary statistics and confidence intervals, standardized differences (rather than *P* values) are a preferred way to exhibit imbalances between groups.

3. *P* values are most meaningful in the context of clear, a priori hypotheses that support the main conclusions of a manuscript.

4. Reporting stand-alone *P* values is discouraged, and preference should be given to presentation and interpretation of effect sizes and their uncertainty (confidence intervals) in the scientific context and in light of other evidence. Crossing a threshold (eg, *P* < .05) by itself constitutes only weak evidence.

5. Researchers should define and interpret effect measures that are clinically relevant. For example, clinical importance is often difficult to establish on the odds ratio scale but is clearer on the risk ratio or absolute risk difference scale.

In summary, following Mark et al,[2] we encourage research-

## Abstract

*P* values and hypothesis testing methods are frequently misused in clinical research. Much of this misuse appears to be owing to the widespread, mistaken belief that they provide simple, reliable, and objective triage tools for separating the true and important from the untrue or unimportant. The primary focus in interpreting therapeutic clinical research data should be on the treatment ("oomph") effect, a metaphorical force that moves patients given an effective treatment to a different clinical state relative to their control counterparts. This effect is assessed using 2 complementary types of statistical measures calculated from the data, namely, effect magnitude or size and precision of the effect size. In a randomized trial, effect size is often summarized using constructs, such as odds ratios, hazard ratios, relative risks, or adverse event rate differences. How large a treatment effect has to be to be consequential is a matter for clinical judgment. The precision of the effect size (conceptually related to the amount of spread in the data) is usually addressed with confidence intervals. *P* values (significance tests) were first proposed as an informal heuristic to help assess how "unexpected" the observed effect size was if the true state of nature was no effect or no difference. Hypothesis testing was a modification of the significance test approach that envisioned controlling the false-positive rate of study results over many (hypothetical) repetitions of the experiment of interest. Both can be helpful but, by themselves, provide only a tunnel vision perspective on study results that ignores the clinical effects the study was conducted to measure.

Section 5

## Dividing Data Comparisons into Categories based on p values

## PROBABLE CAUSE

A P value measures whether an observed result can be attributed to chance. But it cannot answer a researcher's real question: what are the odds that a hypothesis is correct? Those odds depend on how strong the result was and, most importantly, on how plausibile the hypothesis is in the first place.

■ Chance of real effect
■ Chance of no real effect

**THE LONG SHOT**
19-to-1 odds against

**THE TOSS-UP**
1-to-1 odds

**THE GOOD BET**
9-to-1 odds in favour

**Before the experiment**
The plausibility of the hypothesis — the odds of it being true — can be estimated from previous experiments, conjectured mechanisms and other expert knowledge. Three examples are shown here.
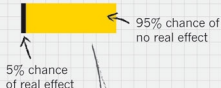
95% chance of no real effect
5% chance of real effect

50%   50%

90%   10%

**The measured P value**
A value of 0.05 is conventionally deemed 'statistically significant'; a value of 0.01 is considered 'very significant'.

$P = 0.05$   $P = 0.01$   $P = 0.05$   $P = 0.01$   $P = 0.05$   $P = 0.01$

11% chance of real effect

**After the experiment**
A small P value can make a hypothesis more plausible, but the difference may not be dramatic.
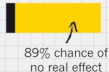
89% chance of no real effect   30%   70%   71%   29%   89%   11%   96%   4%   99%   1%

# Gelman on *p* values, 1

> *The common practice of dividing data comparisons into categories based on significance levels is terrible, but it happens all the time. . . . so it's worth examining the prevalence of this error.*

Consider, for example, this division:

- "really significant" for $p < .01$,
- "significant" for $p < .05$,
- "marginally significant" for $p < .1$, and
- "not at all significant" otherwise.

Now consider some typical *p*-values in these ranges: say, $p = .005$, $p = .03$, $p = .08$, and $p = .2$.

Translate these two-sided *p*-values back into z-scores. . .

*Gelman* 2016-10-15

# Gelman on *p* values, 2

| Description | really sig. | sig. | marginally sig. | not at all sig. |
|---|---|---|---|---|
| *p* value | 0.005 | 0.03 | 0.08 | 0.20 |
| Z score | 2.8 | 2.2 | 1.8 | 1.3 |

The seemingly yawning gap in p-values comparing the not at all significant *p*-value of .2 to the really significant *p*-value of .005, is only a z score of 1.5.

If you had two independent experiments with z-scores of 2.8 and 1.3 and with equal standard errors and you wanted to compare them, you'd get a difference of 1.5 with a standard error of 1.4, which is completely consistent with noise.

# Gelman on *p* values, 3

From a **statistical** point of view, the trouble with using the p-value as a data summary is that the p-value is only interpretable in the context of the null hypothesis of zero effect, and (much of the time), nobody's interested in the null hypothesis.

Indeed, once you see comparisons between large, marginal, and small effects, the null hypothesis is irrelevant, as you want to be comparing effect sizes.

From a **psychological** point of view, the trouble with using the p-value as a data summary is that this is a kind of deterministic thinking, an attempt to convert real uncertainty into firm statements that are just not possible (or, as we would say now, just not replicable).

**The key point**: The difference between statistically significant and NOT statistically significant is not, generally, statistically significant.

*[I]t is unacceptably easy to publish statistically significant evidence consistent with any hypothesis.*

*The culprit is a construct we refer to as **researcher degrees of freedom**. In the course of collecting and analyzing data, researchers have many decisions to make: Should more data be collected? Should some observations be excluded? Which conditions should be combined and which ones compared? Which control variables should be considered? Should specific measures be combined or transformed or both?*

Simmons et al. *link*

# "Researcher Degrees of Freedom", 2

> *. . . It is rare, and sometimes impractical, for researchers to make all these decisions beforehand. Rather, it is common (and accepted practice) for researchers to explore various analytic alternatives, to search for a combination that yields statistical significance, and to then report only what worked. The problem, of course, is that the likelihood of at least one (of many) analyses producing a falsely positive finding at the 5% level is necessarily greater than 5%.*

For more, see

- Gelman's blog $2012 - 11 - 01$ "Researcher Degrees of Freedom",
- Paper by *Simmons* and others, defining the term.

# And this is really hard to deal with. . .

**The garden of forking paths**: Why multiple comparisons can be a problem, even when there is no "fishing expedition" or p-hacking and the research hypothesis was posited ahead of time

> *Researcher degrees of freedom can lead to a multiple comparisons problem, even in settings where researchers perform only a single analysis on their data. The problem is there can be a large number of potential comparisons when the details of data analysis are highly contingent on data, without the researcher having to perform any conscious procedure of fishing or examining multiple p-values. We discuss in the context of several examples of published papers where data-analysis decisions were theoretically-motivated based on previous literature, but where the details of data selection and analysis were not pre-specified and, as a result, were contingent on data.*

- *Link* to the paper from Gelman and Loken

# Abandon Statistical Significance

Gelman blog $2017-09-26$ on "Abandon Statistical Significance"

"Measurement error and variation are concerns even if your estimate is more than 2 standard errors from zero. Indeed, if variation or measurement error are high, then you learn almost nothing from an estimate even if it happens to be 'statistically significant.' "

Read the whole paper *here*

Reviewing "The Association Between Men's Sexist Attitudes and Facial Hair" PubMed 26510427 (*Arch Sex Behavior* May 2016)

Headline Finding: A sample of ~500 men from America and India shows a significant relationship between sexist views and the presence of facial hair.

Excerpt 1:

> *Since a linear relationship has been found between facial hair thickness and perceived masculinity . . . we explored the relationship between facial hair thickness and sexism. . . . Pearson's correlation found no significant relationships between facial hair thickness and hostile or benevolent sexism, education, age, sexual orientation, or relationship status.*

# Facial Hair and Sexist Attitudes

Excerpt 2:

> *We conducted pairwise comparisons between clean-shaven men and each facial hair style on hostile and benevolent sexism scores. . . . For the purpose of further analyses, participants were classified as either clean-shaven or having facial hair based on their self-reported facial hair style . . . There was a significant Facial Hair Status by Sexism Type interaction . . .*

- So their headline finding appeared only because, after their first analysis failed, they shook and shook the data until they found something statistically significant.

# Facial Hair and Sexist Attitudes

Excerpt 2:

> We conducted pairwise comparisons between clean-shaven men
> and each facial hair style on hostile and benevolent sexism scores. .
> . . For the purpose of further analyses, participants were classified
> as either clean-shaven or having facial hair based on their self-
> reported facial hair style . . . There was a significant Facial Hair
> Status by Sexism Type interaction . . .

- So their headline finding appeared only because, after their first
  analysis failed, they shook and shook the data until they found
  something statistically significant.
- All credit to the researchers for admitting that they did this, but poor
  practice of them to present their result in the abstract to their paper
  without making this clear, and too bad that the journal got suckered
  into publishing this.

# How should we react to this?

Gelman:

- Statisticians such as myself should recognize that the point of criticizing a study is, in general, to shed light on statistical errors, maybe with the hope of reforming future statistical education.
- Researchers and policymakers should not just trust what they read in published journals.

## What to do?

In advance, **and** after the fact, think hard about what a plausible effect size might be.

Then. . .

- Analyze *all* your data.
- Present *all* your comparisons, not just a select few.
    - A big table, or even a graph, is what you want.
- Make your data public.
    - If the topic is worth studying, you should want others to be able to make rapid progress.

# A Formula for Decoding Health News

**Health Headlines are Advertising**

Think about this headline: "Hospital checklist cut infections, saved lives."

- Suppose you are a little surprised that a checklist could really save lives. If you think say the odds of this being true are 1 in 4, you would set your initial gut feeling to $1/4$. Because this number is less than one, it means initially you're less likely to believe the study.

**Bayes' Rule**

Final opinion = (initial gut feeling) * (study support for headline)

Jeff Leek, *fivethirtyeight.com*

# Assessing Study Support for a Headline

1. Was the study a clinical study in humans?
2. Was the outcome of the study something directly related to human health like longer life or less disease? Was the outcome something you care about, such as living longer or feeling better?
3. Was the study a randomized, controlled trial (RCT)?
4. Was it a large study - at least hundreds of patients?
5. Did the treatment have a major impact on the outcome?
6. Did predictions hold up in at least two separate groups of people?

## Assessing Study Support

Support for Headline: Multiply by 2 for every yes, and $1/2$ for every no.

# Evaluating A Research Article

Intensive care units (ICUs) at Michigan hospitals implemented a new strategy for reducing infections through training, a daily goals sheet and a program to improve the culture of safety in the ICUs. The doctors measured the rate of infection before and after implementing this safety program.

1. Was the study a clinical study in humans?
   - The study was done in humans in ICUs (+)
2. Was the outcome of the study something directly related to human health like longer life or less disease? Was the outcome something you care about, such as living longer or feeling better?
   - The outcome was the rate of infections after surgery. According to the article, these infections cost U.S. hospitals up to $2.3 billion annually. (+)

## Evaluating a Research Article

3. Was the study a randomized, controlled trial (RCT)?
   - The study compared the same hospitals before and after a change in ICU policy. This is an example of a crossover study, which is not as strong as a randomized trial but does account for some of the differences among hospitals because the same ICUs were measured before and after using the checklist. (-)
4. Was it a large study - at least hundreds of patients?
   - The study looked at more than 100 ICUs over 1,981 months. In total, it followed patients for 375,757 catheter-days. (A catheter-day means watching one patient for one day while she is on a catheter.) This is a huge number of days to monitor patients for potential infections. (+)
5. Did the treatment have a major impact on the outcome?
   - The study showed a sustained drop of up to 66 percent in infections. (+)
6. Did predictions hold up in at least two separate groups of people?
   - The study looked at 103 hospitals in Michigan. (+)

So we have 5 + and 1 - in our evaluation of this study.

# Final Opinion?

- So, a large study showed a major drop in infections that is directly medically important. For the sake of the exercise, let's multiply by two every time we see a *yes* answer and by $1/2$ every time we see a *no* answer. I would say this study's result is about 16 times more likely (five out of six *yes* answers so 2 x 2 x 2 x 2 x 2 x (1/2) = 16) if checklists really do reduce infections than if they don't. I set study support for headline = 16.

- Multiply to get final opinion on headline = 1/4*16 = 4, also expressed as 4/1. I would say that my updated odds are 4 to 1 that the headline is true. The strength of the study won over my initially skeptical gut feeling.

## Bayes' Rule

Final opinion = (initial gut feeling) * (study support for headline)
Source: Jeff Leek, *fivethirtyeight.com*

# Evaluating Health News: For Consumers

1. Watch out for single source stories. They're usually based on a press release, which will have a hidden agenda.
2. Beware of stories that don't mention cost. It's crucial information. (If the cost of the great, new treatment is out of reach – it's not that great, is it?)
3. Headline percentages are misleading. If something "reduces your risk of X by 50%" chances are that number doesn't mean what you think it means.
4. If it sounds too good to be true, it probably is. If a report presents only or primarily the benefits of a new treatment, it's a bad report. ALL healthcare interventions have trade-offs.
5. Patient anecdotes are not data. Beware of stories that rely on them. Anecdotes are used to compensate for data that are unavailable or flawed.

Source: *NPR*

## Evaluating Health News: For Consumers

6. A "simple screening test" is never simple. The decision to take one is one of the most complex and difficult decisions a health consumer can make.

7. Watch out for hyperbolic language. "Breakthrough", "first-of-its-kind", and "game-changer" are red flags. When you read "it may become..." substitute "it may not become..."

8. Newer isn't always better. Often the latest test, treatment or procedure is no better than what already exists, just pricier.

9. Beware of disease-mongering. Risk factors, symptoms for diseases, or data can be exaggerated in a way that causes needless worry, and expense.

10. The latest treatment may not exist yet, or ever. "Awaiting FDA approval" or "in pre-clinical trial phase" means it's still a pipe dream.

11. There is a leap from mice to men. Getting from rodent trials to human use is a very, very long road, that may in fact lead nowhere.

Source: NPR

Section 6

**Being A More Effective / Transparent / Reproducible / Open Source Scientist**

From *PLoS Comput Biol* *link*

EDITORIAL

# Ten Simple Rules for Effective Statistical Practice

Robert E. Kass[1], Brian S. Caffo[2], Marie Davidian[3], Xiao-Li Meng[4], Bin Yu[5], Nancy Reid[6]*

## Rule 10: Make Your Analysis Reproducible

# Goals of Reproducible Research

The goal of reproducible research is to tie specific instructions to data analysis so that scholarship can be recreated, better understood and verified. This is usually facilitated by literate programming - a document that combines content and data analytic code. Software? R and R Studio, mostly:

1. Be able to reproduce your own results and allow others to reproduce your results
2. Reproduce an entire report / manuscript / thesis / book / website with a single system command when changes occur (in operating system, statistical software, graphics engines, source data, derived variables, analysis, interpretation).
3. Save time.
4. Provide the ultimate documentation of work done.

Vanderbilt *Tutorial*

Karl -- this is very interesting, however you used an old version of the data (n=143 rather than n=226).

I'm really sorry you did all that work on the incomplete dataset.

Bruce

# Five Practical Tips

1. Document everything.
2. Everything is a (text) file.
3. All files should be human-readable.
4. Explicitly tie your files together.
5. Have a plan to organize, store and make your files available.

The papers/slideshows/abstracts are not the research. The research is the full software environment, code and data that produced the results. (Donoho, 2010). Separating research from its advertisement makes it hard for others to verify or reproduce our findings.

Your closest collaborator is you, six months ago, but you don't respond to emails. (Holder via Broman)

Karl Broman, Steps Towards Reproducible Research

# Wisdom from DL Donoho (2010) re: Open-Source

But other people will use my data and code to compete with me?

- True.

# Wisdom from DL Donoho (2010) re: Open-Source

But other people will use my data and code to compete with me?

- True.
- But competition means that strangers will read your work, try to learn from you, cite you, and try to do things even better.
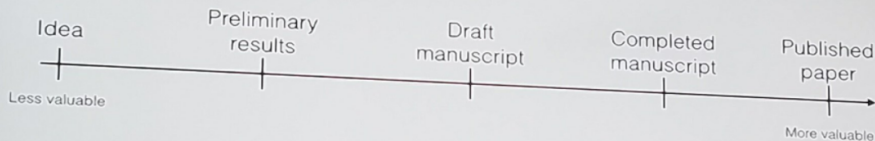
# Wisdom from DL Donoho (2010) re: Open-Source

But other people will use my data and code to compete with me?

- True.
- But competition means that strangers will read your work, try to learn from you, cite you, and try to do things even better.
- If you prefer obscurity, why are you publishing?

How I thought of my goals in grad school:

Idea — Preliminary results — Draft manuscript — Completed manuscript — Published paper

Less valuable → More valuable

How I should have been thinking of them:

Anything still on your computer
(Data, code, results, draft, finished paper)

Anything out in the world
(Paper, preprint, product, blog post, open source, tweet)

Less valuable → More valuable

# What I Think I Think Now

- Null hypothesis significance testing is much harder than I thought.
  - The null hypothesis is almost never a real thing.
  - Rather than rejiggering the cutoff, I would mostly abandon the *p* value as a summary
  - Replication is far more useful than I thought it was.
- Some hills aren't worth dying on.
  - Think about uncertainty intervals more than confidence or credible intervals
  - Retrospective calculations about Type S (sign) and Type M (magnitude) errors can help me illustrate ideas.
- Which method to use is far less important than finding better data
  - The biggest mistake I make regularly is throwing away useful data
  - I'm not the only one with this problem.
- The best thing I do most days is communicate more clearly.
  - When stuck in a design, I think about how to get better data.
  - When stuck in an analysis, I try to turn a table into a graph.
- I have A LOT to learn.

# Be The Change You Want To See In The World

- Today's slides are at https://github.com/THOMASELOVE/rethink.
- You'll find all of the references there, as well, and some other sources.

## Thomas E. Love, Ph.D.

- Professor of Medicine, Population & Quantitative Health Sciences, CWRU School of Medicine
- Director of Biostatistics and Data Science, Population Health Research Institute, The MetroHealth System
- Chief Data Scientist, Better Health Partnership
- Fellow, American Statistical Association

My email is Thomas dot Love at case dot edu.