

# NEJM Statistical Guidelines for Authors

Under the Hood

Dave Harrington

18 September 2019

# MY COORDINATES

- Department of Biostatistics, Harvard T.H. Chan School of Public Health
- Department of Data Sciences, Dana-Farber Cancer Institute
- Statistical Consultant, New England Journal of Medicine
- [davidharrington@g.harvard.edu](mailto:davidharrington@g.harvard.edu)

# GUIDELINES WERE JOINT EFFORT

## Statistical Consultants

- Ralph D'Agostino, BU
- Constantine Gatsonis, Brown
- D Harrington, HSPH, DFCI
- Joe Hogan, Brown
- David Hunter, Oxford
- Sharon-Lise Normand, HMS

Jeff Drazen, former Editor-in-Chief Mary Beth Hamel, Executive Deputy Editor

All Deputy Editors read and approved guidelines

# STATISTICAL REVIEW AT NEJM

Six paid statistical consultants (reviewers)

We meet weekly with the editors

Review all research articles that editors decide to move through the system

- We see approximately 5% of the 5,000 - 6000 submitted articles

Statistical reviewers

- Help set statistical 'policy'
- Work toward consistency in reviews

# BACKGROUND

Revised author guidelines on statistical reporting posted on NEJM website July 1, 2019

<https://www.nejm.org/author-center/new-manuscripts>

Editorial describing the guidelines published 18 July 2019

- N Engl J Med 2019; 381:285-286

Guidelines cover many aspects of statistical reporting

Editorial emphasized section on p-values

# OUR GOALS

Emphasize aspects of good practice in analysis and reporting

- All points discussed have been in the statistical literature for decades

Add consistency to

- Statistical reviewing at NEJM
- Standards of analysis in reporting of results

Move toward standards of reporting in observational studies

- Elements of STROBE (Strengthening of Reporting of Observational Studies in Epidemiology)
- <https://strobe-statement.org/>

# GOALS . . .

Show the data

- Point estimates and confidence intervals
- Reduce irrelevant  $p$ -values

Acknowledge that we cannot influence design

Every good paper should have a path to publication, even with

- No planned adjustment for multiplicity
- No specified method for missing data

# GOALS . . .

Minimize misleading or misunderstood statistical jargon

- Nominal p-value
- Nominally significant
- 'Controlling' for confounders in regression models



# OUR CONSTRAINTS

For busy clinicians, conclusions in a paper should be clear, direct, and supportable.

Statistical background of NEJM readers

Word limits in NEJM make it difficult to discuss nuance.

- 250 words for abstracts
- 2700 words for article

Normal review time and weekly publication schedule make transition to major changes difficult.

# CONTENT OF THE GUIDELINES

For all studies

- Describe sample size and power calculations
- $p$ -values should be two-sided.

For RCTs:

- Provide protocol and statistical analysis plan (SAP).
- Analysis should match design
- Specify and follow multiplicity adjustments in SAP
- No plan for multiplicity adjustment: report point estimates with confidence intervals, no  $p$ -values

# THE GUIDELINES . . .

RCTs:

- Forest plots for subgroups, but no  $p$ -values
- OK to use unadjusted  $p$ -values for safety comparisons
- Report event rates as well as hazard ratios
- Provide CONSORT diagram

# THE GUIDELINES

## Observational studies:

- Provide a SAP if you have one
- No SAP - provide a clear analysis plan in methods section or supplement
- No plan for multiplicity adjustment: report point estimates with confidence intervals, no  $p$ -values
- Show distributions of the measured confounders
- Model diagnostics and sensitivity analyses in the supplement.
- If possible, retest findings in a validation dataset.

# WHAT ABOUT THOSE UNADJUSTED CONFIDENCE INTERVALS?

Is a confidence interval always a test in disguise?

Interpretation of confidence intervals by readers

Typical insert I recommend to authors for methods:

*Because the statistical analysis plan did not include a provision for correcting for multiplicity, secondary and other outcomes are reported as point estimates and 95% confidence intervals. The widths of the confidence intervals have not been adjusted for multiplicity, so the intervals should not be used to infer definitive treatment effects for secondary outcomes.*

# POST HOC ADJUSTMENTS FOR MULTIPLICITY

Personal communication from Rebecca Betensky

For each of 10,000 replicates in a simulation

- Simulate collection of correlated  $p$ -values
- Choose a method of FWER or FDR control that is most favorable to the analyst
  - Minimize the minimum  $p$ -value or
  - Maximize the number of adjusted  $p$ -values  $< 0.05$
- Average FWER or FDR across replicates

Choice of methods from

- FWER: Bonferroni, Holm, Hochberg, Hommel, Šidák single step, Šidák step down.
- FDR: Benjamini-Hochberg, Benjamini-Yekutieli, two-stage Benjamini-Hochberg.

# RESULTS

number of tests	correlation	maximize number of adjusted $p < 0.05$			minimize minimum adjusted $p$		
		procedures used			procedures used		
		FWER	FDR	ALL	FWER	FDR	ALL
20	0	0.052	0.052	0.053	0.003	0.050	0.053
	0.5	0.032	0.037	0.037	0.001	0.035	0.037
	0.9	0.019	0.026	0.026	0.000	0.026	0.026
50	0	0.046	0.046	0.047	0.009	0.039	0.047
	0.5	0.029	0.034	0.034	0.004	0.031	0.034
	0.9	0.015	0.025	0.025	0.000	0.025	0.025
100	0	0.051	0.052	0.053	0.028	0.025	0.053
	0.5	0.027	0.034	0.034	0.008	0.026	0.034
	0.9	0.010	0.022	0.022	0.000	0.022	0.022

8.

# PERILS OF POST HOC ADJUSTMENT

Does not distinguish between confirmation and discovery

Are all the primary secondary and subgroup analyses part of the family of tests examining a treatment effect?

- If so, reduces power for primary outcome



# UNADJUSTED $p$ -VALUES IN SAFETY OUTCOMES

From Metra, et al. N Engl J Med 2019; 381:716-726 DOI: 10.1056/NEJMoa1801291

**Table 3.** Most Frequent Adverse Events up to and Including Day 5.\*

Event	Serelaxin Group (N = 3257)	Placebo Group (N = 3248)	P Value†
	no. (%)		
Hypokalemia	263 (8.1)	242 (7.5)	0.35
Cardiac failure	192 (5.9)	215 (6.6)	0.23
Muscle spasms	80 (2.5)	49 (1.5)	0.006
Hypotension	77 (2.4)	65 (2.0)	0.32
Headache	74 (2.3)	92 (2.8)	0.15
Constipation	70 (2.1)	57 (1.8)	0.25
Urinary tract infection	63 (1.9)	68 (2.1)	0.65

\* Shown are events that occurred in at least 2% of patients in either group.

† P values were calculated with the use of the Cochran–Mantel–Haenszel test.

## $p$ VALUES FOR INTERACTIONS

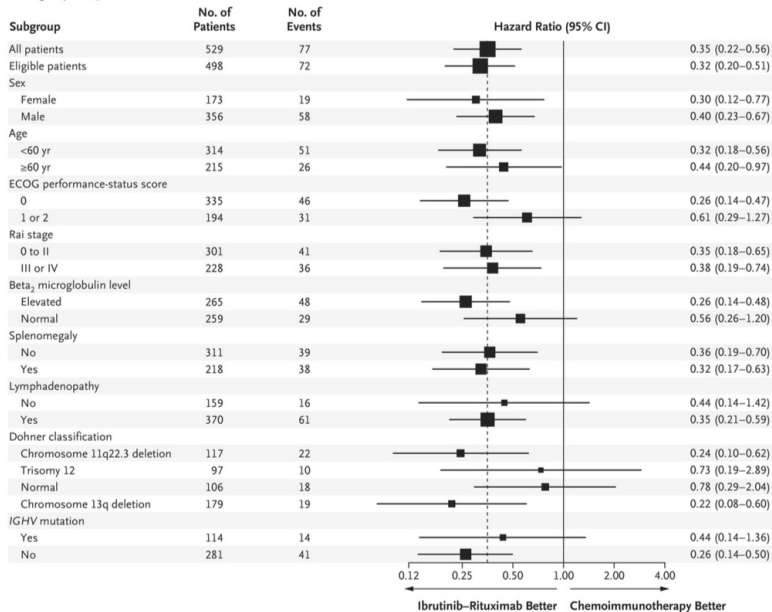
Did we do more harm than good in eliminating  $p$ -values in forest plots of subgroups?

Deputy editors, authors and some statisticians reluctant to eliminate these

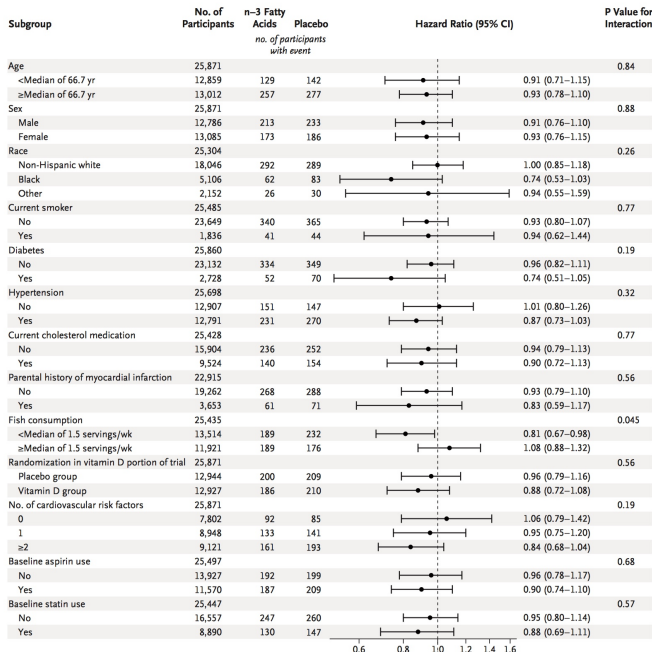
- Authors and editors found lists of non-significant  $p$ -values reassuring.
- Easy to recommend that significant  $p$ -values be ignored.
- No easy general recommendation for point estimates and confidence intervals for interaction effects.

# NEW FORM OF FOREST PLOTS

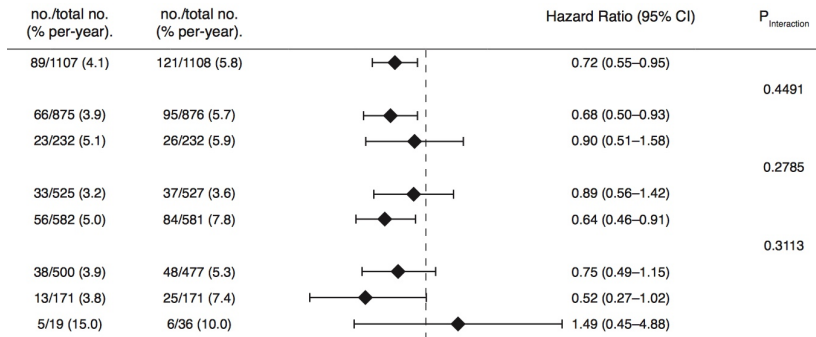
## C Subgroup Analysis



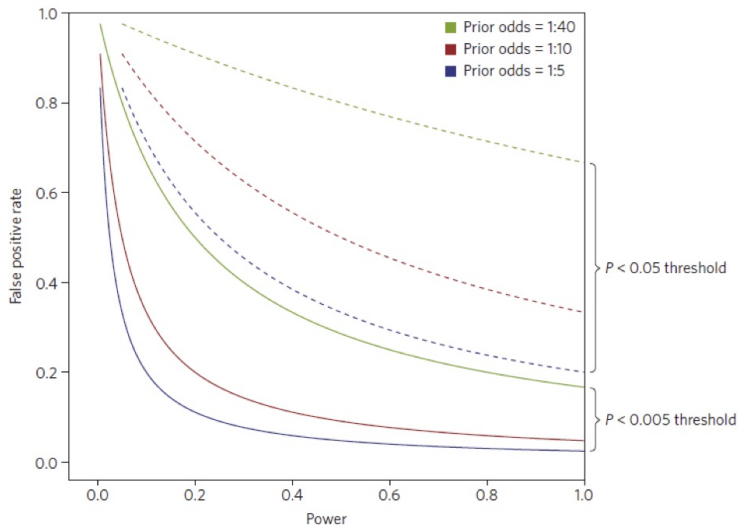
# BEFORE CHANGE



# PERHAPS P-VALUES ARE HELPFUL



# CALCULATING POSTERIOR ODDS (NAT. HUMAN BEHAVIOR, JAN 2018)



## EFFECT OF INCREASING $\alpha$

Prior odds = 1:10, power = 0.80

Alpha	False Pos. Prob
0.10	0.56
0.15	0.65
0.20	0.71
0.25	0.76
0.40	0.83
0.60	0.88

False Pos. Prob = probability of incorrectly claiming alternative is true, given data

# ASSUMPTIONS BEHIND THE CALCULATIONS

Decision to use an intervention based on single trial, single p-value

Prior odds is knowable, at least approximately

There is a simple alternative for which power is relevant

For practicing clinician, is reproducibility the right metric?



# HAVE WE MADE UNSUPPORTED CLAIMS EASIER?

Common to encounter phrases in drafts such as

*A consistent pattern for improved survival . . . was noted across multiple subgroups*

or

*An exploratory . . . analysis . . . of patients who resumed [intervention] was superior to that of patients who did not*

PDF of slides available at under mskcc\_2019 at  
<https://github.com/dave-harrington/talks>

# READING

Benjamin DJ, et al. Redefine statistical significance. Nat Hum Behavior 2018;2:6-10. doi: 10.1038/s41562-017-0189-z

Ioannidis JPA. Retiring statistical significance would give bias a free pass. Nature. 2019;567 (7749):461. doi:10.1038/d41586-019-00969-2

National Academies of Sciences, Engineering, and Medicine. Reproducibility and replicability in science. Washington, DC: National Academies Press, 2019. <http://nap.edu/25303>

Wasserstein RL, Schirm AL, Lazar NA. Moving to a world beyond  $p < 0.05$ . Am Stat. 2019;73:1-19. doi:10.1080/00031305.2019.1583913

Dmitrienko A, Bretz F, Westfall PH, et al. Multiple testing methodology. In: Dmitrienko A, Tamhane AC, Bretz F, eds. Multiple testing problems in pharmaceutical statistics. New York: Chapman and Hall/CRC Press, 2009:35-98.