

The Need for Equivalence Testing in Economics

Jack Fitzgerald

Vrije Universiteit Amsterdam

September 13, 2024



“No Detectable Effects”

Bessone et al. (2021, QJE): Sleep improvement RCT with 400 people in Chennai, India

“No Detectable Effects”

Bessone et al. (2021, QJE): Sleep improvement RCT with 400 people in Chennai, India

- | At baseline, avg. participant has sleep patterns mirroring clinical insomnia

“No Detectable Effects”

Bessone et al. (2021, QJE): Sleep improvement RCT with 400 people in Chennai, India

- | At baseline, avg. participant has sleep patterns mirroring clinical insomnia
- | The intervention is very effective (27 extra minutes of night sleep)

“No Detectable Effects”

Bessone et al. (2021, QJE): Sleep improvement RCT with 400 people in Chennai, India

- | At baseline, avg. participant has sleep patterns mirroring clinical insomnia
- | The intervention is very effective (27 extra minutes of night sleep)

However, per their abstract...

*\Contrary to expert predictions and a large body of sleep research, increased nighttime sleep had **no detectable effects** on cognition, productivity, decision making, or well being... ”*

“No Detectable Effects”

Bessone et al. (2021, QJE): Sleep improvement RCT with 400 people in Chennai, India

- | At baseline, avg. participant has sleep patterns mirroring clinical insomnia
- | The intervention is very effective (27 extra minutes of night sleep)

However, per their abstract...

*\Contrary to expert predictions and a large body of sleep research, increased nighttime sleep had **no detectable effects** on cognition, productivity, decision making, or well being... ”*

By their own admission, these findings contradict expert priors and large bodies of research

“No Detectable Effects”

Bessone et al. (2021, QJE): Sleep improvement RCT with 400 people in Chennai, India

- | At baseline, avg. participant has sleep patterns mirroring clinical insomnia
- | The intervention is very effective (27 extra minutes of night sleep)

However, per their abstract...

*“Contrary to expert predictions and a large body of sleep research, increased nighttime sleep had **no detectable effects** on cognition, productivity, decision making, or well being...”*

By their own admission, these findings contradict expert priors and large bodies of research

- | So what do they mean by ‘**no detectable effects**?’

Null Estimates in Bessone et al. (2021)

What they mean: Results are not stat. sig. different from zero

Null Estimates in Bessone et al. (2021)

What they mean: Results are not stat. sig. different from zero

- | **They are not alone** in interpreting insignificant results in this way

This Happens All the Time

From 2020-2023, 279 null claims made in abstracts of 158 articles in T5 journals are defended by statistically insignificant results [Detailed Results](#)

This Happens All the Time

From 2020-2023, 279 null claims made in abstracts of 158 articles in T5 journals are defended by statistically insignificant results [Detailed Results](#)

- | > 72% of these null claims aren't qualified by references to statistical significance, estimate magnitudes, or a lack of evidence

This Happens All the Time

From 2020-2023, 279 null claims made in abstracts of 158 articles in T5 journals are defended by statistically insignificant results [Detailed Results](#)

- | > 72% of these null claims aren't qualified by references to statistical significance, estimate magnitudes, or a lack of evidence

Researchers and readers interpret such findings as evidence of null/negligible relationships (McShane & Gal 2016, McShane & Gal 2017)

Why Is This a Problem?

Generally inferring that stat. insig. results are null results is known to be bad scientific practice (Altman & Bland 1995; Imai, King, & Stuart 2008; Wasserstein & Lazar 2016)

- I Statistical insignificance may just reflect imprecision

Why Is This a Problem?

Generally inferring that stat. insig. results are null results is known to be bad scientific practice (Altman & Bland 1995; Imai, King, & Stuart 2008; Wasserstein & Lazar 2016)

- | Statistical insignificance may just reflect imprecision

Under standard NHST, null results and imprecision are conflated. Credibility problems follow:

Why Is This a Problem?

Generally inferring that stat. insig. results are null results is known to be bad scientific practice (Altman & Bland 1995; Imai, King, & Stuart 2008; Wasserstein & Lazar 2016)

- | Statistical insignificance may just reflect imprecision

Under standard NHST, null results and imprecision are conflated. Credibility problems follow:

- | Null result penalty from beliefs of low quality and unpublishability (McShane & Gal 2016; McShane & Gal 2017; Chopra et al. 2024)

Why Is This a Problem?

Generally inferring that stat. insig. results are null results is known to be bad scientific practice (Altman & Bland 1995; Imai, King, & Stuart 2008; Wasserstein & Lazar 2016)

- | Statistical insignificance may just reflect imprecision

Under standard NHST, null results and imprecision are conflated. Credibility problems follow:

- | Null result penalty from beliefs of low quality and unpublishability (McShane & Gal 2016; McShane & Gal 2017; Chopra et al. 2024)
- | Publication bias from non-publication of null results (Fanelli 2012; Franco, Malhotra, & Simonovits 2014; Andrews & Kasy 2019)

Why Is This a Problem?

Generally inferring that stat. insig. results are null results is known to be bad scientific practice (Altman & Bland 1995; Imai, King, & Stuart 2008; Wasserstein & Lazar 2016)

- | Statistical insignificance may just reflect imprecision

Under standard NHST, null results and imprecision are conflated. Credibility problems follow:

- | Null result penalty from beliefs of low quality and unpublishability (McShane & Gal 2016; McShane & Gal 2017; Chopra et al. 2024)
- | Publication bias from non-publication of null results (Fanelli 2012; Franco, Malhotra, & Simonovits 2014; Andrews & Kasy 2019)
- | High Type II error rates, given current practices and power levels (Ioannidis, Stanley, & Doucouliagos 2017; Askarov et al. 2023)

Why Is This a Problem?

Generally inferring that stat. insig. results are null results is known to be bad scientific practice (Altman & Bland 1995; Imai, King, & Stuart 2008; Wasserstein & Lazar 2016)

- | Statistical insignificance may just reflect imprecision

Under standard NHST, null results and imprecision are conflated. Credibility problems follow:

- | Null result penalty from beliefs of low quality and unpublishability (McShane & Gal 2016; McShane & Gal 2017; Chopra et al. 2024)
- | Publication bias from non-publication of null results (Fanelli 2012; Franco, Malhotra, & Simonovits 2014; Andrews & Kasy 2019)
- | High Type II error rates, given current practices and power levels (Ioannidis, Stanley, & Doucouliagos 2017; Askarov et al. 2023)

It doesn't have to be this way.

Equivalence Testing in a Nutshell

1. Set a region around zero wherein relationship of interest should be practically equivalent to zero (i.e., economically insignificant)

Equivalence Testing in a Nutshell

1. Set a region around zero wherein relationship of interest should be practically equivalent to zero (i.e., economically insignificant)
2. Use interval tests to assess if sig. bounded within this region

Equivalence Testing in a Nutshell

1. Set a region around zero wherein relationship of interest should be practically equivalent to zero (i.e., economically insignificant)
2. Use interval tests to assess if sig. bounded within this region

Common in medicine, political science, and psychology (see e.g., Piaggio et al. 2012; Hartman & Hidalgo 2018; Lakens, Scheel, & Isager 2018)

This Project

What is equivalence testing?

- I introduce simple frequentist equivalence testing techniques to economists

This Project

What is equivalence testing?

- | I introduce simple frequentist equivalence testing techniques to economists

Why do we need to use it?

- | 36-63% of estimates defending null claims in top economics journals fail lenient equivalence tests
- | Type II error rates in economics are likely quite high

This Project

What is equivalence testing?

- | I introduce simple frequentist equivalence testing techniques to economists

Why do we need to use it?

- | 36-63% of estimates defending null claims in top economics journals fail lenient equivalence tests
- | Type II error rates in economics are likely quite high

How do we perform equivalence testing credibly?

- | I develop software commands and guidelines for credible and relatively easy implementation

The Wrong Hypotheses: NHST

Standard NHST hypotheses :

$$H_0 : = 0$$

$$H_A : \neq 0$$

The Wrong Hypotheses: NHST

Standard NHST hypotheses :

$$H_0 : = 0$$

$$H_A : \neq 0$$

When trying to show that $= 0$ using NHST, two key problems:

The Wrong Hypotheses: NHST

Standard NHST hypotheses :

$$H_0 : = 0$$

$$H_A : \neq 0$$

When trying to show that $= 0$ using NHST, two key problems:

1. The burden of proof is shifted : Researchers start by assuming they're right

The Wrong Hypotheses: NHST

Standard NHST hypotheses :

$$H_0 : = 0$$

$$H_A : \neq 0$$

When trying to show that $= 0$ using NHST, two key problems:

1. The burden of proof is shifted : Researchers start by assuming they're right
2. Imprecision is `good': Less precision → higher chance of stat. insig. results

The Wrong Hypotheses: NHST

Standard NHST hypotheses :

$$H_0 : = 0$$

$$H_A : \neq 0$$

When trying to show that $= 0$ using NHST, two key problems:

1. The burden of proof is shifted : Researchers start by assuming they're right
2. Imprecision is 'good': Less precision → higher chance of stat. insig. results

It's thus a logical fallacy to generally infer that stat. insig. results are null results (appeal to ignorance)

The Right Hypotheses: Equivalence Testing

We'll fix these problems by 1) flipping the hypotheses and 2) relaxing the constraints.
As a reminder, NHST hypotheses:

$$H_0 : \mu = 0$$

$$H_A : \mu \neq 0$$

And now equivalence testing hypotheses:

$$H_0 : \mu \in 0$$

$$H_A : \mu \notin 0$$

The Right Hypotheses: Equivalence Testing

We'll fix these problems by 1) flipping the hypotheses and 2) relaxing the constraints.
As a reminder, NHST hypotheses:

$$H_0 : \mu = 0$$

$$H_A : \mu \neq 0$$

And now equivalence testing hypotheses:

$$H_0 : \mu \notin 0$$

$$H_A : \mu = 0$$

If we can set a range of values $[-\epsilon, +\epsilon]$ wherein $\mu = 0$, then we can find stat. sig. evidence for H_A with a simple interval test

The Equivalence Testing Framework

We begin by setting a range of values $[-, +]$, where $- < +$, called the region of practical equivalence (ROPE)

The Equivalence Testing Framework

We begin by setting a range of values $[-\epsilon, \epsilon]$, where $\epsilon < \delta$, called the region of practical equivalence (ROPE)

- ! The ROPE is the range of values we'd call economically insignificant

The Equivalence Testing Framework

We begin by setting a range of values $[-, +]$, where $- < +$, called the region of practical equivalence (ROPE)

- | The ROPE is the range of values we'd call economically insignificant
- | This is a subjective judgment call that will differ for different relationships of interest
- | I show how to credibly aggregate ROPEs later in this talk [Credible ROPE-Setting](#)

The Equivalence Testing Framework

We begin by setting a range of values $[\delta, \epsilon]$, where $\delta < \epsilon$, called the region of practical equivalence (ROPE)

- ! The ROPE is the range of values we'd call economically insignificant
- ! This is a subjective judgment call that will differ for different relationships of interest
- ! I show how to credibly aggregate ROPEs later in this talk [Credible ROPE-Setting](#)

Once we have a ROPE, we can set up the equivalence testing hypotheses:

$$H_0 : \theta \in [\delta, \epsilon]$$

$$H_A : \theta \notin [\delta, \epsilon]$$

Two One-Sided Tests (TOST)

We can identically write the equivalence testing hypotheses as

$$H_0 : < \quad \text{or} \quad > \quad +$$

$$H_A : \quad \text{and} \quad +$$

Two One-Sided Tests (TOST)

We can identically write the equivalence testing hypotheses as

$$H_0 : < \quad \text{or} \quad > \quad +$$

$$H_A : \quad \text{and} \quad +$$

Further, we can assess the joint H_A using two one-sided tests:

$$H_0 : <$$

$$H_A :$$

$$H_0 : > \quad +$$

$$H_A : \quad +$$

Two One-Sided Tests (TOST)

We can identically write the equivalence testing hypotheses as

$$H_0 : < \quad \text{or} \quad > \quad +$$

$$H_A : \quad \text{and} \quad +$$

Further, we can assess the joint H_A using two one-sided tests:

$$H_0 : <$$

$$H_0 : > +$$

$$H_A :$$

$$H_A : +$$

Stat. sig. evidence for both H_A statements using one-sided tests is stat. sig. evidence that 0
(Schuirmann 1987; Berger & Hsu 1996, [Procedural Details](#) [Visualization](#))

Equivalence Confidence Intervals (ECIs)

\hat{s}^2 (1 - α) equivalence confidence interval (ECI) is just its (1 - 2α) CI

- | If \hat{s}^2 (1 - α) ECI is entirely bounded in the ROPE, then we have size- α evidence under the TOST procedure that $\theta = 0$ (Berger & Hsu 1996)

Revisiting Bessone et al. (2021)

Estimates defending null claims should be significantly bounded within reasonably wide ROPEs

Revisiting Bessone et al. (2021)

Estimates defending null claims should be significantly bounded within reasonably wide ROPEs

- | However, 28% of the 'null' estimates in Bessone et al. (2021) aren't significantly bounded beneath $j = 0:2$
- | 71% aren't significantly bounded beneath $j = 0:1$

Revisiting Bessone et al. (2021)

Estimates defending null claims should be significantly bounded within reasonably wide ROPEs

- | However, 28% of the 'null' estimates in Bessone et al. (2021) aren't significantly bounded beneath $j = 0:2$
- | 71% aren't significantly bounded beneath $j = 0:1$

Takeaway: Bessone et al. (2021) cannot guarantee precise nulls for a large proportion of their 'null' estimates, which 'fail' lenient equivalence tests

Data

1. Systematically-selected replication sample

- | 876 estimates defending 135 null claims in abstracts of 81 articles in T5 economics journals published from 2020-2021 [Claim Example](#)
- | Estimates defending these null claims are reproducible with publicly-available data

Data

1. Systematically-selected replication sample

- | 876 estimates defending 135 null claims in abstracts of 81 articles in T5 economics journals published from 2020-2021 [Claim Example](#)
- | Estimates defending these null claims are reproducible with publicly-available data

2. Prediction platform data

- | I survey 62 researchers on the Social Science Prediction Platform for predictions and judgments on equivalence testing results in my sample

Equivalence Testing Failure Rates

I compute avg.equivalence testing failure rates
in the replication sample

Equivalence Testing Failure Rates

I compute avg. equivalence testing failure rates in the replication sample

- I First ROPE: $r \sim N(0, 1; 0, 1)$
- I $|r| = 0, 1$ is larger than over 25% of published results in economics (Doucouliagos 2011)

Effect Size Standardization

Equivalence Testing Failure Rates

I compute avg. equivalence testing failure rates in the replication sample

- I First ROPE: $r \sim 2 [0:1; 0:1]$
- I $|r_j| = 0:1$ is larger than over 25% of published results in economics (Doucouliagos 2011)
Effect Size Standardization
- I Second ROPE: $r \sim 2 [0:2; 0:2]$
- I $|r_j| = 0:2$ is quite large for economic effect sizes
Benchmarking Sample

Equivalence Testing Failure Rates

I compute avg. equivalence testing failure rates in the replication sample

- I First ROPE: $r \sim 2 [0:1; 0:1]$
- I $|r_j| = 0:1$ is larger than over 25% of published results in economics (Doucouliagos 2011)
Effect Size Standardization
- I Second ROPE: $r \sim 2 [0:2; 0:2]$
- I $|r_j| = 0:2$ is quite large for economic effect sizes
Benchmarking Sample

Models defending null claims in T5 journals should have no trouble significantly bounding estimates within ROPEs this wide

Equivalence Testing Failure Rates are Unacceptably High

Equivalence testing failure rates range from 36-63%. [Robustness Checks](#) [TST Framework](#) [Mechanisms](#)

Equivalence Testing Failure Rates are Unacceptably High

Equivalence testing failure rates range from 36-63%. [Robustness Checks](#) [TST Framework](#) [Mechanisms](#)

- Interpretation : 62% of estimates defending the average null claim can't significantly bound their estimates beneath $\tau_j = 0:1$ (see Model 4)

Failure Curves

Equivalence testing failure rates stay unacceptably high even as ROPEs become ridiculously large

Failure Curves

- Equivalence testing failure rates stay unacceptably high even as ROPEs become ridiculously large
- ! To obtain acceptable failure rates, you'd need to argue that $0.317r_j$ is practically equal to zero

Failure Curves

Equivalence testing failure rates stay unacceptably high even as ROPEs become ridiculously large

- | To obtain acceptable failure rates, you'd need to argue that θ_0 is practically equal to zero
- | θ_0 is larger than nearly 75% of published effects in economics (Doucouliagos 2011)

Researchers Anticipate Unacceptably High Failure Rates

The median researcher finds failure rates from 11-13% acceptable, but (pretty accurately) predicts failure rates from 35-38%. Takeaways:

Researchers Anticipate Unacceptably High Failure Rates

The median researcher finds failure rates from 11-13% acceptable, but (pretty accurately) predicts failure rates from 35-38%. Takeaways:

1. Researchers don't trust null results under standard NHST, but this mistrust is well-placed

Researchers Anticipate Unacceptably High Failure Rates

The median researcher finds failure rates from 11-13% acceptable, but (pretty accurately) predicts failure rates from 35-38%. Takeaways:

1. Researchers don't trust null results under standard NHST, but this mistrust is well-placed
2. More credible testing frameworks are necessary to restore trust

Credible ROPE-Setting

ROPEs need to be set independently to be credible (Lange & Freitag 2005; Ofori et al. 2023)

- | `ROPE-hacking' is a key concern

Credible ROPE-Setting

ROPEs need to be set independently to be credible (Lange & Freitag 2005; Ofori et al. 2023)

- | `ROPE-hacking' is a key concern
- | To maintain independence & credibility, you shouldn't set your ROPEs { you should get other people to set them for you

Credible ROPE-Setting

ROPEs need to be set independently to be credible (Lange & Freitag 2005; Ofori et al. 2023)

- | `ROPE-hacking' is a key concern
- | To maintain independence & credibility, you shouldn't set your ROPEs { you should get other people to set them for you

Solution: Survey independent experts/stakeholders for their judgments

Credible ROPE-Setting

ROPEs need to be set independently to be credible (Lange & Freitag 2005; Ofori et al. 2023)

- | `ROPE-hacking' is a key concern
- | To maintain independence & credibility, you shouldn't set your ROPEs { you should get other people to set them for you

Solution: Survey independent experts/stakeholders for their judgments

- | Practically feasible using online platforms such as the Social Science Prediction Platform (DellaVigna, Pope, & Vivaldi 2019)

Credible ROPE-Setting

ROPEs need to be set independently to be credible (Lange & Freitag 2005; Ofori et al. 2023)

- | `ROPE-hacking' is a key concern
- | To maintain independence & credibility, you shouldn't set your ROPEs { you should get other people to set them for you

Solution: Survey independent experts/stakeholders for their judgments

- | Practically feasible using online platforms such as the Social Science Prediction Platform (DellaVigna, Pope, & Vivaldi 2019)
- | Example from this project : Alongside predictions of failure rates, I elicit what failure rates researchers deem acceptable

The Equivalence Testing Framework

Software Commands & More Information

tsti Stata command (QR:
Github)

equivtest R package, containing
tst command (QR: Github)

Working paper (QR: PDF link,
personal website)

Website: <https://jack-tzgerald.github.io>
Email: j.f.tzgerald@vu.nl

References I

Altman, D. G. and J. M. Bland (1995).

Statistics notes: Absence of evidence is not evidence of absence.

BMJ 311(7003), 485{485.

Andrews, I. and M. Kasy (2019).

Identification of and correction for publication bias.

American Economic Review 10(8), 2766{2794.

Askarov, Z., A. Doucouliagos, H. Doucouliagos, and T. D. Stanley (2023).

Selective and (mis)leading economics journals: Meta-research evidence.

Journal of Economic Surveys, Forthcoming

Berger, R. L. and J. C. Hsu (1996).

Bioequivalence trials, intersection-union tests and equivalence confidence sets.

Statistical Science 1(4).

References II

Chopra, F., I. Haaland, C. Roth, and A. Stegmann (2024).

The null result penalty.

The Economic Journal 134(657), 193{219.

Cohen, J. (1988).

Statistical power analysis for the behavioral sciences (2 ed.).

L. Erlbaum Associates.

DellaVigna, S., D. Pope, and E. Vivaldi (2019).

Predict science to improve science.

Science 366(6464), 428{429.

Doucoulagos, H. (2011).

How large is large? Preliminary and relative guidelines for interpreting partial correlations in economics.

Working Paper SWP 2011/5, Deakin University, Geelong, Australia.

References III

Dreber, A., M. Johannesson, and Y. Yang (2024).

Selective reporting of placebo tests in top economics journals.

Economic Inquiry.

Fanelli, D. (2012).

Negative results are disappearing from most disciplines and countries.

Scientometrics 90(3), 891{904.

Finner, H. and K. Strassburger (2002).

The partitioning principle: A powerful tool in multiple decision theory.

The Annals of Statistics 30(4), 1194{1213.

Goeman, J. J., A. Solari, and T. Stijnen (2010).

Three-sided hypothesis testing: Simultaneous testing of superiority, equivalence and inferiority.

Statistics in Medicine 29(20), 2117{2125.

References IV

Hartman, E. and F. D. Hidalgo (2018).

An equivalence approach to balance and placebo tests.

American Journal of Political Science 6(4), 1000{1013.

Imai, K., G. King, and E. A. Stuart (2008).

Misunderstandings between experimentalists and observationalists about causal inference.

Journal of the Royal Statistical Society Series A: Statistics in Society 17(2), 481{502.

Ioannidis, J. P., T. D. Stanley, and H. Doucouliagos (2017).

The power of bias in economics research.

The Economic Journal 127(605).

Lakens, D., A. M. Scheel, and P. M. Isager (2018).

Equivalence testing for psychological research: A tutorial.

Advances in Methods and Practices in Psychological Science 2(2), 259{269.

References V

Lange, S. and G. Freitag (2005).

Choice of delta: Requirements and reality { results of a systematic review.
Biometrical Journal 47(1), 12{27.

McShane, B. B. and D. Gal (2016).

Blinding us to the obvious? The effect of statistical training on the evaluation of evidence.
Management Science 62(8), 1707{1718.

McShane, B. B. and D. Gal (2017).

Statistical significance and the dichotomization of evidence.
Journal of the American Statistical Association 112(519), 885{895.

References VI

Ofori, S., T. Cafaro, P. Devereaux, M. Marcucci, L. Mbuagbaw, L. Thabane, and G. Guyatt (2023).

Noninferiority margins exceed superiority effect estimates for mortality in cardiovascular trials in high-impact journals.

Journal of Clinical Epidemiology 16,120{27.

Piaggio, G., D. R. Elbourne, S. J. Pocock, S. J. Evans, and D. G. Altman (2012).

Reporting of noninferiority and equivalence randomized trials.

JAMA 308(24), 2594{2604.

Schuirman, D. J. (1987).

A comparison of the two one-sided tests procedure and the power approach for assessing the equivalence of average bioavailability.

Journal of Pharmacokinetics and Biopharmaceutics 1(6), 657{680.

References VII

Sha er, J. P. (1986).

Modi ed sequentially rejective multiple test procedures.

Journal of the American Statistical Association 8(395), 826{831.

Wasserstein, R. L. and N. A. Lazar (2016).

The ASA statement on p-values: Context, process, and purpose.

The American Statistician 70(2), 129{133.

Null Claim Classification

This Happens All the Time

The TOST Procedure

First, compute test statistics

$$t = \frac{\hat{\mu}}{s}$$

$$t_+ = \frac{\hat{\mu}_+}{s}$$

The relevant test statistic is the smaller of the two:

$$t_{\text{TOST}} = \arg \min_{t, t_+} t, t_+$$

The critical value for a size- α TOST procedure is the one-sided critical value t_{α}

1. If $t_{\text{TOST}} = t$, then there is stat. sig. evidence that $\mu \in [\mu_-, \mu_+]$ if $t \geq t_{\alpha}$
2. If $t_{\text{TOST}} = t_+$, then there is stat. sig. evidence that $\mu \in [\mu_-, \mu_+]$ if $t_+ \geq t_{\alpha}$

A single TOST procedure maintains size α even without multiple hypothesis corrections (Berger & Hsu 1996)

TOST Concept

TOST Example

Claim Example

The bolded text represents the two null claims made by this abstract:

\This article estimates peer effects originating from the ability composition of tutorial groups for undergraduate students in economics. We manipulated the composition of groups to achieve a wide range of support, and assigned students conditional on their prior ability randomly to these groups. The data support a specification in which the impact of group composition on achievement is captured by the mean and standard deviation of peers' prior ability, their interaction, and interactions with students' own prior ability. When we assess the aggregate implications of these peer effects regressions for group assignment, we find that low- and medium-ability students gain on an average 0.19 SD units of achievement by switching from ability mixing to three-way tracking. Their dropout rate is reduced by 12 percentage points (relative to a mean of 0.6). High-ability students are unaffected. Analysis of survey data indicates that in tracked groups, low-ability students have more positive interactions with other students, and are more involved. We find no evidence that teachers adjust their teaching to the composition of groups. "

Data

Standardized Effect Sizes

I aggregate all regression results into two effect size measures

1. Standardized coefficients :

$$s = \begin{cases} \frac{\beta}{s_Y} & \text{if } D \text{ is binary} \\ \frac{\beta_D}{s_Y} & \text{otherwise} \end{cases} \quad S = \begin{cases} \frac{SE(\beta)}{s_Y} & \text{if } D \text{ is binary} \\ \frac{SE(\beta)_D}{s_Y} & \text{otherwise} \end{cases}$$

s_Y and s_D are respectively within-sample SDs of Y and D

S is closely related to the classical Cohen's effect size

2. Partial correlation coefficients (PCCs) :

$$r = \rho \frac{t_{NHST}}{t_{NHST}^2 + df} \quad SE(r) = \frac{1}{\rho} \frac{r^2}{df}$$

t_{NHST} is the usual t -statistic and df is degrees of freedom

r PCCs are widely-used in economic meta-analyses

Benchmarking Sample

The Three-Sided Testing Framework

The three-sided testing (TST) framework (Goeman, Solari, & Stijnen 2010) uses ROPE [; +] to assess 's practical significance using three tests:

The Three-Sided Testing Framework

The three-sided testing (TST) framework (Goeman, Solari, & Stijnen 2010) uses ROPE [; +] to assess θ 's practical significance using three tests:

1. Two-sided test: Is $\theta < -\delta$?
2. TOST procedure: Is $\theta \in [-\delta; +\delta]$?
3. Two-sided test: Is $\theta > +\delta$?

The Three-Sided Testing Framework

The three-sided testing (TST) framework (Goeman, Solari, & Stijnen 2010) uses ROPE [$;$ $+$] to assess θ 's practical significance using three tests:

1. Two-sided test: $|\theta| < \tau$?
2. TOST procedure: $|\theta| \leq \tau$ [$;$ $+$]?
3. Two-sided test: $|\theta| > \tau$?

Significance conclusions can be derived from the smallest of these three p -values (Shafer 1986; Finner & Strassburger 2002)

The Three-Sided Testing Framework

The three-sided testing (TST) framework (Goeman, Solari, & Stijnen 2010) uses ROPE [$;$ $+$] to assess θ 's practical significance using three tests:

1. Two-sided test: Is $|\theta| < \delta$?
2. TOST procedure: Is $\theta \in [L; +]$?
3. Two-sided test: Is $|\theta| > \delta + \epsilon$?

Significance conclusions can be derived from the smallest of these three p-values (Shafer 1986; Finner & Strassburger 2002)

- 1. If no p-value $< \alpha$, then results are inconclusive the researcher must stay agnostic about the practical significance of θ

The Three-Sided Testing Framework

The **three-sided testing (TST) framework** (Goeman, Solari, & Stijnen 2010) uses ROPE [; +] to assess 's practical significance using three tests:

1. Two-sided test: Is $\theta < \tau$?
2. TOST procedure: Is $\theta \geq [\tau^- ; +]$?
3. Two-sided test: Is $\theta > \tau^+$?

Significance conclusions can be derived from the smallest of these three p -values (Shafer 1986; Finner & Strassburger 2002)

- | If no p -value $< \alpha$, then results are *inconclusive*: the researcher must stay agnostic about the practical significance of
- | Embracing this uncertainty may be uncomfortable/limiting, but my results show that standard practice tolerates high error rates

The Three-Sided Testing Framework

The **three-sided testing (TST) framework** (Goeman, Solari, & Stijnen 2010) uses ROPE [μ ; $\mu + \delta$] to assess μ 's practical significance using three tests:

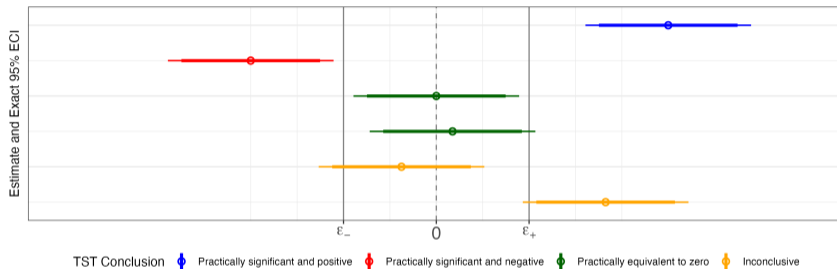
1. Two-sided test: Is $\mu < \mu - \delta$?
2. TOST procedure: Is $\mu \geq [\mu - \delta; \mu + \delta]$?
3. Two-sided test: Is $\mu > \mu + \delta$?

Significance conclusions can be derived from the smallest of these three p -values (Shafer 1986; Finner & Strassburger 2002)

- | If no p -value $< \alpha$, then results are *inconclusive*: the researcher must stay agnostic about the practical significance of μ
- | Embracing this uncertainty may be uncomfortable/limiting, but my results show that standard practice tolerates high error rates

Example from this project: I show that my failure rates are significantly bounded above the median failure rates that researchers deem acceptable [Main Results](#)

The TST Framework Visualized



Under TST, given their 95% ECIs and CIs, these estimates are respectively:

- | Practically significant and above the ROPE
- | Practically significant and below the ROPE
- | Practically equivalent to zero
- | Inconclusive

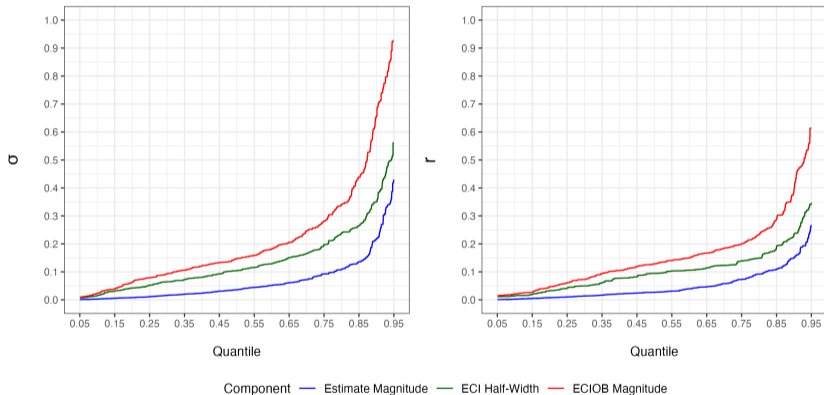
Failure Rate Robustness

These failure rates remain large and significant when...

- | Switching from r to r
- | Switching from exact to asymptotically approximate tests
- | Switching aggregation procedures
- | Removing initially stat. sig. estimates
- | Separating models by regressor type combination (i.e., binary vs. non-binary)
- | Removing non-replicable estimates from the sample
- | Removing models that require conformability modifications from the sample (e.g., logit/probit models put through `margins`, `dydx()`)

Main Results

Mechanisms



Power is a greater driver of equivalence testing failure rates than effect size

Main Results