The Need for Equivalence Testing in Economics

Jack Fitzgerald

Vrije Universiteit Amsterdam

November 25, 2024



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

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

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

Bessone et al. (2021, QJE): Sleep improvement RCT with \approx 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)

Bessone et al. (2021, QJE): Sleep improvement RCT with pprox 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..."

Bessone et al. (2021, QJE): Sleep improvement RCT with pprox 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

Bessone et al. (2021, QJE): Sleep improvement RCT with pprox 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?'

What: Equivalence Testi 000000 Why: Results

Null Estimates in Bessone et al. (2021)



Significant NHST Result? - Yes - No

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

What: Equivalence Testi

Why: Results

Null Estimates in Bessone et al. (2021)



Significant NHST Result? - Yes - No

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

Abstract

Smallholder farming in many developing countries is characterized by low productivity and lowquality output. Low quality limits the price farmers an command and their potential income. We conduct a series of experiments among maize farmers in Uganda to shed light on the barriers to quality upgrading and to study its potential. We find that the causal return to quality is zero. Providing access to a market where quality is paid a market premium led to an increase in farm productivity and income from farming. Our findings reveal the importance of demand-side constraints in limiting rural income and productivity growth.

Abstract

Consumers rely on the price changes of goods in their grocery bundles when forming expectations about aggregate inflation. We use micro data that uniquely match individual expectations, detailed information about consumption bundles, and item-level prices. The weights consumers assign to price changes dopend on the frequency of purchase, rather than expenditure share, and positive price changes loom larger than eaglity price changes. Prices of good officed in the same store but not purchased do not affect inflation expectations, nor do other dimensions. Our results provide empirical guidance for models of expectations formation with hereogneous consumers.

Abstract

We study how political turnover in mayoral elections in Brazil affects public service provision by local governments. Exploiting a regression discontinuity design for close elections, we find that municipalities with a new party in office experience upheavals in the municipal bureaucracy: new persionel are appointed across multiple service sectors, and at both managerial and non-managerial levels. In education, the increase in the replacement rate of personnel in a hools controlled by the municipal government is accompanied by test scores that are 0.05–0.08 standard deviations lower. In contrast, turnover of the mayor's party does not impact local frommunicipal) schools. These findings suggest that political turnover can adversely affect the quality of public services when the bureaucracy is not shielded from the political process.

Abstract

This paper estimates intertemporal labor supply responses to two-year long income tax holidays staggered across Swiss cantons. Cantons shifted from an income tax system based on the previous two years' income to a standard annual pay as you earn system, leaving two years of income untaxed. We find significant but quantitatively very small responses of wage earnings with an intertemporal elasticity of 0.025 overall. High wage income earners and especially the self-employed display larger responses with elasticities around 0.1 and 0.25, respectively, most likely driven by tax avoldance. We find no effects along the extensive margin at all.

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

Abstract

Smallholder farming in many developing countries is characterized by low productivity and lowquality output. Low quality limits the price farmers an command and their potential income. We conduct a series of experiments among maize farmers in Uganda to shed light on the barriers to quality upgrading and to study its potential. We find that the causal return to quality is zero. Providing access to a market where quality is paid a market premium led to an increase in farm productivity and income from farming. Our findings reveal the importance of demand-side constraints in limiting rural income and productivity growth.

Abstract

Consumers rely on the price changes of goods in their grocery bundles when forming expectations about aggregate inflation. We use micro data that uniquely match individual expectations, detailed information about consumption bundles, and item-level prices. The weights consumers assign to price changes dopend on the frequency of purchase, rather than expenditure share, and positive price changes loom larger than eaglity price changes. Prices of good officed in the same store but not purchased do not affect inflation expectations, nor do other dimensions. Our results provide empirical guidance for models of expectations formation with hereogneous consumers.

Abstract

We study how political turnover in mayoral elections in Brazil affects public service provision by local governments. Exploiting a regression discontinuity design for close elections, we find that municipalities with a new party in office experience upheavals in the municipal bureaucray: new personnel are appointed across multiple service sectors, and at both managerial and non-managerial levels. In education, the increase in the replacement rate of personnel in a chools controlled by the municipal government is accompanied by test scores that are 0.05–0.08 standard deviations lower. In contrast, turnover of the mayor's party does not impact local Inonmunicipal) schools. These findings suggest that political turnover can adversely affect the quality of public services when the bureaucrays is not shielded from the political process.

Abstract

This paper estimates intertemporal labor supply responses to two-year long income tax holidays staggered across Swiss cantons. Cantons shifted from an income tax system based on the previous two years' income to a standard annual pay as you earn system, leaving two years of income untaxed. We find significant but quantitatively very small responses of wage earnings with an intertemporal elasticity of 0.025 overall. High wage income earners and especially the self-employed display larger responses with elasticities around 0.1 and 0.25, respectively, most likely driven by tax avoldance. We find no effects along the extensive margin at all.

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

Abstract

Smallholder farming in many developing countries is characterized by low productivity and lowquality output. Low quality limits the price farmers an command and their potential income. We conduct a series of experiments among maize farmers in Uganda to shed light on the barriers to quality upgrading and to study its potential. We find that the causal return to quality is zero. Providing access to a market where quality is paid a market premium led to an increase in farm productivity and income from farming. Our findings reveal the importance of demand-side constraints in limiting rural income and productivity growth.

Abstract

Consumers rely on the price changes of goods in their grocery bundles when forming expectations about aggregate inflation. We use micro data that uniquely match individual expectations, detailed information about consumption bundles, and item-level prices. The weights consumers assign to price changes dopend on the frequency of purchase, rather than expenditure share, and positive price changes loom larger than eaglity price changes. Prices of good officed in the same store but not purchased do not affect inflation expectations, nor do other dimensions. Our results provide empirical guidance for models of expectations formation with hereogneous consumers.

Abstract

We study how political turnover in mayoral elections in Brazil affects public service provision by local governments. Exploiting a regression discontinuity design for close elections, we find that municipalities with a new party in office experience upheavals in the municipal bureaucray: new personnel are appointed across multiple service sectors, and at both managerial and non-managerial levels. In education, the increase in the replacement rate of personnel in a chools controlled by the municipal government is accompanied by test scores that are 0.05–0.08 standard deviations lower. In contrast, turnover of the mayor's party does not impact local Inonmunicipal) schools. These findings suggest that political turnover can adversely affect the quality of public services when the bureaucrays is not shielded from the political process.

Abstract

This paper estimates intertemporal labor supply responses to two-year long income tax holidays staggered across Swiss cantons. Cantons shifted from an income tax system based on the previous two years' income to a standard annual pay as you earn system, leaving two years of income untaxed. We find significan but quantitatively very small responses of wage earnings with an intertemporal elasticity of 0.025 overall. High wage income earners and especially the self-employed display larger responses with elasticities around 0.1 and 0.25, respectively, most likely driven by tax avoldance. We find no effects along the extensive margin at all.

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)

Vrije Universiteit Amsterdam

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

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 δ would be **practically** equivalent to zero (i.e., *economically insignificant*)

Equivalence Testing in a Nutshell



- 1. Set a region around zero wherein relationship of interest δ would be **practically** equivalent to zero (i.e., *economically insignificant*)
- 2. Use interval tests to assess if $\hat{\delta}$ is sig. bounded within this region

Equivalence Testing in a Nutshell



- 1. Set a region around zero wherein relationship of interest δ would be **practically** equivalent to zero (i.e., *economically insignificant*)
- 2. Use interval tests to assess if $\hat{\delta}$ is 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:

$$H_0: \delta = 0$$
$$H_A: \delta \neq 0$$

The Wrong Hypotheses: NHST

Standard NHST hypotheses:

$$H_0: \delta = 0$$
$$H_A: \delta \neq 0$$

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

The Wrong Hypotheses: NHST

Standard NHST hypotheses:

$$H_0: \delta = 0$$
$$H_A: \delta \neq 0$$

When trying to show that $\delta = 0$ 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:

$$H_0: \delta = 0$$
$$H_A: \delta \neq 0$$

When trying to show that $\delta = 0$ 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 \rightarrow higher chance of stat. insig. results

The Wrong Hypotheses: NHST

Standard NHST hypotheses:

$$H_0: \delta = 0$$
$$H_A: \delta \neq 0$$

When trying to show that $\delta = 0$ 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 \rightarrow 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 contraints. As a reminder, **NHST hypotheses**:

 $H_0: \delta = 0$ $H_A: \delta \neq 0$

And now equivalence testing hypotheses:

 $H_0: \delta \not\approx 0$ $H_A: \delta \approx 0$

The Right Hypotheses: Equivalence Testing

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

 $H_0: \delta = 0$ $H_A: \delta \neq 0$

And now equivalence testing hypotheses:

 $H_0: \delta \not\approx 0$ $H_A: \delta \approx 0$

If we can set a range of values $[\epsilon_-, \epsilon_+]$ wherein $\delta \approx 0$, 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 $[\epsilon_-, \epsilon_+]$, where $\epsilon_- < \epsilon_+$, called the *region of practical equivalence (ROPE)*

The Equivalence Testing Framework

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

• The ROPE is the range of δ values we'd call *economically insignificant*

The Equivalence Testing Framework

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

- The ROPE is the range of δ 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 $[\epsilon_-, \epsilon_+]$, where $\epsilon_- < \epsilon_+$, called the *region of practical* equivalence (ROPE)

- The ROPE is the range of δ 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:

 $H_0: \delta \notin [\epsilon_-, \epsilon_+]$ $H_A: \delta \in [\epsilon_-, \epsilon_+]$

Two One-Sided Tests (TOST)



We can identically write the equivalence testing hypotheses as

 $H_0: \delta < \epsilon_- \text{ or } \delta > \epsilon_+$ $H_A: \delta \ge \epsilon_- \text{ and } \delta \le \epsilon_+$
Two One-Sided Tests (TOST)



We can identically write the equivalence testing hypotheses as

 $H_0: \delta < \epsilon_- \text{ or } \delta > \epsilon_+$ $H_A: \delta \ge \epsilon_- \text{ and } \delta \le \epsilon_+$

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

$$\begin{array}{ll} H_0: \, \delta < \epsilon_- & H_0: \, \delta > \epsilon_+ \\ H_A: \, \delta \ge \epsilon_- & H_A: \, \delta \le \epsilon_+ \end{array}$$

Two One-Sided Tests (TOST)



We can identically write the equivalence testing hypotheses as

 $H_0: \delta < \epsilon_- \text{ or } \delta > \epsilon_+$ $H_A: \delta \ge \epsilon_- \text{ and } \delta \le \epsilon_+$

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

 $\begin{array}{ll} H_0: \, \delta < \epsilon_- & H_0: \, \delta > \epsilon_+ \\ H_A: \, \delta \geq \epsilon_- & H_A: \, \delta \leq \epsilon_+ \end{array}$

Stat. sig. evidence for **both** H_A statements using one-sided tests is stat. sig. evidence that $\delta \approx 0$ (Schuirmann 1987; Berger & Hsu 1996) Procedural Details Visualization

Why: Results

Equivalence Confidence Intervals (ECIs)



 $\hat{\delta}$'s $(1-\alpha)$ equivalence confidence interval (ECI) is just its $(1-2\alpha)$ Cl

► If $\hat{\delta}$'s $(1 - \alpha)$ ECI is entirely bounded in the ROPE, then we have size- α evidence under the TOST procedure that $\delta \approx 0$ (Berger & Hsu 1996) Comparison w/ TOST

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 |σ| = 0.2
- 71% aren't significantly bounded beneath |r| = 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 |σ| = 0.2
- ▶ 71% aren't significantly bounded beneath |r| = 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 2020-2023 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 2020-2023 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

- ▶ First ROPE: $r \in [-0.1, 0.1]$
- ► |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

- ▶ First ROPE: $r \in [-0.1, 0.1]$
- ► |r| = 0.1 is larger than over 25% of published results in economics (Doucouliagos 2011) Effect Size Standardization
- **Second ROPE**: $\sigma \in [-0.2, 0.2]$
- ► |σ| = 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

- ▶ First ROPE: $r \in [-0.1, 0.1]$
- ► |r| = 0.1 is larger than over 25% of published results in economics (Doucouliagos 2011) Effect Size Standardization
- Second ROPE: $\sigma \in [-0.2, 0.2]$
- |σ| = 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

Many 'Null' Estimates Fail Lenient Equivalence Tests



Over 39% of the 'null' estimates in my sample can't be significantly bounded beneath 0.2 σ

Many 'Null' Estimates Fail Lenient Equivalence Tests



Over 39% of the 'null' estimates in my sample can't be significantly bounded beneath 0.2 σ

• Over 69% can't be significantly bounded beneath 0.1r

Equivalence Testing Failure Rates are Unacceptably High



Equivalence testing failure rates range from 36-63% Robustness Checks TST Framework Mechanis

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 |r| = 0.1 (see Model 4)

Why: Results

Failure Curves



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

Why: Results 000000

Failure Curves



Effect Size Measure 📕 σ 📃 r

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| is practically equal to zero

Why: Results

Failure Curves



Effect Size Measure 📒 σ 📃 r

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| is practically equal to zero
- ▶ |0.317*r*| 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

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

► 'ROPE-hacking' is a key concern

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

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

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, & Vivalt 2019)

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, & Vivalt 2019)
- **Example from this project**: Alongside predictions of failure rates, I elicit what failure rates researchers deem acceptable

The Equivalence Testing Framework

The Next Step: Practical Significance Testing

Natural to want to combine equivalence testing with tests for δ 's practical significance

The Next Step: Practical Significance Testing

Natural to want to combine equivalence testing with tests for $\delta{}^\prime{\rm s}$ practical significance

► Can be done using the three-sided testing (TST) framework (Goeman, Solari, & Stijnen 2010)

The Next Step: Practical Significance Testing

Natural to want to combine equivalence testing with tests for δ 's practical significance

► Can be done using the three-sided testing (TST) framework (Goeman, Solari, & Stijnen 2010)

Given ROPE $[\epsilon_-, \epsilon_+]$, the idea is to assess δ 's practical significance using three tests:

The Next Step: Practical Significance Testing

Natural to want to combine equivalence testing with tests for $\delta '{\rm s}$ practical significance

► Can be done using the three-sided testing (TST) framework (Goeman, Solari, & Stijnen 2010)

Given ROPE $[\epsilon_-, \epsilon_+]$, the idea is to assess δ 's practical significance using three tests:

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

The Next Step: Practical Significance Testing

Natural to want to combine equivalence testing with tests for $\delta '{\rm s}$ practical significance

► Can be done using the three-sided testing (TST) framework (Goeman, Solari, & Stijnen 2010)

Given ROPE $[\epsilon_-, \epsilon_+]$, the idea is to assess δ 's practical significance using three tests:

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

Significance conclusions can be derived from the smallest of these three *p*-values

If no p-value < α, then results are *inconclusive*: the researcher must stay agnostic about the practical significance of δ

The Next Step: Practical Significance Testing

Natural to want to combine equivalence testing with tests for $\delta '{\rm s}$ practical significance

► Can be done using the three-sided testing (TST) framework (Goeman, Solari, & Stijnen 2010)

Given ROPE $[\epsilon_-, \epsilon_+]$, the idea is to assess δ 's practical significance using three tests:

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

Significance conclusions can be derived from the smallest of these three *p*-values

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

The Next Step: Practical Significance Testing

Natural to want to combine equivalence testing with tests for $\delta '{\rm s}$ practical significance

► Can be done using the three-sided testing (TST) framework (Goeman, Solari, & Stijnen 2010)

Given ROPE $[\epsilon_-, \epsilon_+]$, the idea is to assess δ 's practical significance using three tests:

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

Significance conclusions can be derived from the smallest of these three *p*-values

- If no *p*-value $< \alpha$, then results are *inconclusive*: the researcher must stay agnostic about the practical significance of δ
- 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 Three-Sided Testing 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

Main Results

ShinyTST App



Peder Isager & I are currently working on a TST tutorial, complete with a Shiny app Main Results



ShinyTST App
What: Equivalence Testi 000000 Why: Results

How: The Future

Software Commands & More Information





tsti Stata command (QR: Github) eqtesting R package, containing tst command (QR: Github)

Working paper (QR: PDF link, personal website)

Website: https://jack-fitzgerald.github.io Email: j.f.fitzgerald@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 109(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 11*(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.).



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

Predict science to improve science.

Science 366(6464), 428-429.

Doucouliagos, 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*.

ㅣ 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

1	_	
1		
1		1
1		
1		

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

An equivalence approach to balance and placebo tests. American Journal of Political Science 62(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 171(2), 481–502.

loannidis, 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 1(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(6), 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 161, 20–27.

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.

Schuirmann, 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 15(6), 657–680.

References VII



Shaffer, J. P. (1986).

Modified sequentially rejective multiple test procedures.

Journal of the American Statistical Association 81(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

Category	Claim Type	Example	# Claims	% of Claims
1	Claim that a relationship/phenomenon does not exist or is negligible	D has no effect on Y .	111	39.8%
2	Claim that a relationship/phenomenon does not exist or is negligible, qualified by reference to statistical significance	${\cal D}$ has no significant effect on $Y.$	33	11.8%
3	Claim that a relationship/phenomenon does not exist or is negligible, qualified by reference to something other than statistical significance	${\cal D}$ has no meaningful effect on $Y.$	24	8.6%
4	Claim that a relationship/phenomenon does not (meaningfully) hold in a given direction	${\cal D}$ has no positive effect on $Y.$	53	19%
5	Claim that a relationship/phenomenon does not (meaningfully) hold in a given direction, qualified by reference to statistical significance	${\cal D}$ has no significant positive effect on $Y.$	4	1.4%
6	Claim that a relationship/phenomenon does not (meaningfully) hold in a given direction, qualified by reference to something other than statistical significance	${\cal D}$ has no meaningful positive effect on $Y.$	5	1.8%
7	Claim that there is a lack of evidence for a (meaningful) relationship/phenomenon	There is no evidence that D has an effect on $\boldsymbol{Y}.$	10	3.6%
8	Claim that a variable holds similar values regardless of the values of another variable	\boldsymbol{Y} is similar for those in the treatment group and the control group.	7	2.5%
9	Claim that a relationship/phenomenon holds only or primarily in a subset of the data	The effect of D on Y is concentrated in older respondents.	22	7.9%
10	Claim that a relationship/phenomenon stabilizes for some values of another variable	D has a short term effect on Y that dissipates after Z months.	10	3.6%
	Unqualified null claim	Categories 1, 4, or 8-10	203	72.8%
	Qualified null claim	Categories 2-3 or 5-7	76	27.2%

This Happens All the Time

Vrije Universiteit Amsterdam

The TOST Procedure

First, compute test statistics

$$t_{-}=rac{\hat{\delta}-\epsilon_{-}}{s} \qquad \qquad t_{+}=rac{\hat{\delta}-\epsilon_{+}}{s}$$

The relevant test statistic is the smaller of the two:

$$t_{\mathsf{TOST}} = \argmin_{t \in \{t_-, t_+\}} \{|t|\}$$

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 $\delta \in [\epsilon_{-}, \epsilon_{+}]$ iff $t_{-} \geq t_{\alpha}^{*}$
- 2. If $t_{\text{TOST}} = t_+$, then there is stat. sig. evidence that $\delta \in [\epsilon_-, \epsilon_+]$ iff $t_+ \leq -t_{\alpha}^*$

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:

$$\sigma = \begin{cases} \frac{\delta}{\sigma_Y} \text{ if } D \text{ is binary} \\ \frac{\delta\sigma_D}{\sigma_Y} \text{ otherwise} \end{cases} \qquad \qquad s = \begin{cases} \frac{\mathsf{SE}(\delta)}{\sigma_Y} \text{ if } D \text{ is binary} \\ \frac{\mathsf{SE}(\delta)\sigma_D}{\sigma_Y} \text{ otherwise} \end{cases}$$

 σ_{Y} and σ_{D} are respectively within-sample SDs of Y and D

 $\blacktriangleright \ \sigma$ is closely related to the classical Cohen's d effect size

2. Partial correlation coefficients (PCCs):

$$r = rac{t_{
m NHST}}{\sqrt{t_{
m NHST}^2 + df}}$$
 ${
m SE}(r) = rac{1-r^2}{\sqrt{df}}.$

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

PCCs are widely-used in economic meta-analyses

Failure Rates Introduction

Appendix 00000000

Benchmarking Sample

Article	Setting	Outcome Variable	Exposure Variable	Initial p-Value	σ	r	Location
Acemoglu & Restrepo (2020)	Difference-in-differences analysis of U.S. commuting zones, 1990-2007	Employment rates (continuous)	Industrial robot exposure (continuous)	0.000	-0.206	-0.16	Table 7, Panel A, US exposure to robots, Model 3
Acemoglu et al. (2019)	Difference-in-differences analysis of countries, 1960-2010	Short-run log GDP levels (continuous)	Democratization (binary)	0.001	0.005	0.255	Table 2, Democracy, Model 3
Berman et al. (2017)	African 0.5×0.5 longitude-latitude cells with mineral mines, 1997-2010	Conflict incidence (binary)	Log price of main mineral (continuous)	0.012	0.521	0.007	Table 2, ln price x mines > 0 , Model 1
Deschènes, Greenstone, & Shapiro (2017)	Difference-in-differences analysis of U.S. counties, 2001-2007	Nitrogen dioxide emissions (continuous)	Nitrogen dioxide cap-and-trade participation (binary)	0.000	-0.134	-0.468	Table 2, Panel A, NOx, Model 3
Haushofer & Shaprio (2016)	Experiment with low-income Kenyan households, 2011-2013	Non-durable consumption (continuous)	Unconditional cash transfer (binary)	0.000	0.376	0.195	Table V, Non-durable expenditure, Model 1
Benhassine et al. (2015)	Experiment with families of Moroccan primary school-aged students, 2008-2010	School attendance (binary)	Educational cash transfer to fathers (binary)	0.000	0.18	0.252	Table 5, Panel A, Attending school by end of year 2, among those 6-15 at baseline, Impact of LCT to fathers
Bloom et al. (2015)	Field experiment with Chinese workers, 2010-2011	Attrition (binary)	Voluntarily working from home (binary)	0.002	-0.397	-0.196	Table VIII, Treatment, Model 1
Duflo, Dupas, & Kremer (2015)	Experiment with Kenyan primary school-aged girls, 2003-2010	Reaching eighth grade (binary)	Education subsidy (binary)	0.023	0.1	0.125	Table 3, Panel A, Stand-alone education subsidy, Model 1
Hanushek et al. (2015)	OECD adult workers, 2011-2012	Log hourly wages (continuous)	Numeracy skills (continuous)	0.000	0.091	0.316	Table 5, Numeracy, Model 1
Oswald, Proto, & Sgroi (2015)	UK students, piece-rate laboratory task	Productivity (continuous)	Happiness (continuous)	0.018	0.753	0.244	Table 2, Change in happiness, Model 4



Vrije Universiteit Amsterdam

Failure Rate Robustness

These failure rates remain large and significant when...

- Switching from σ 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

Vrije Universiteit Amsterdam

Jack Fitzgerald