



# The journey toward Population-level Effect Estimation

Martijn Schuemie, PhD  
Janssen Research and Development

---



# Population-level effect estimation

- What is the effect of treatment A on outcome X?
- What is the effect of treatment A on outcome X, compared to exposure B?



# Population-level effect estimation

Evidence  
Generation

- How to produce evidence from the data?

Evidence  
Evaluation

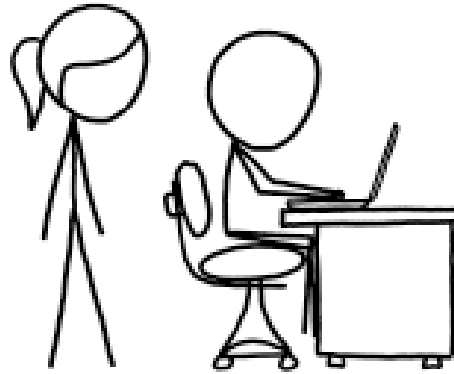
- How do we know the evidence is reliable?

Evidence  
Dissemination

- How do we share evidence to inform decision making?

Doctor, I'm starting on duloxetine,  
should I be worried about stroke?

Let me see what I find in the literature...





# Evidence from literature

Paper by Lee et al, 2016

- Compare new users of SNRIs (includes duloxetine) vs SSRIs
- Taiwanese insurance claims data
- 12 month washout
- remove people using both drugs
- remove people with a prior history of head injury
- remove people with a prior history of stroke or intracranial hemorrhage
- Propensity score: logistic regression with treatment as dependent variable
- HOI is Stroke: first hospitalization with ICD-9 433,434, or 436
- time-varying Cox regression using 5 PS strata

	Crude Hazard Ratio (95% CI)	<i>P</i>	Adjusted Hazard Ratio <sup>a</sup> (95% CI)	<i>P</i>
<b>Main analyses</b>				
SNRIs (n = 76,920) vs SSRIs (n = 582,650)				
Ischemic stroke	0.92 (0.83–1.02)	.12	1.01 (0.90–1.12)	.91



# How reliable is this evidence?

- Can the results be reproduced?
- Did the analysis program do what it was supposed to do?
- Is the estimate unbiased?
- Does the p-value have nominal characteristics?
- Does the confidence interval really represent the uncertainty about the effect size?

Are we really 95% confident the true effect size is between 0.90 and 1.12?



# Population-level effect estimation



- How to produce evidence from the data?



# 'Replicating' Lee et al.

Our replication:

- Compare new users of Duloxetine (SNRI) vs. Sertraline (SSRI)
- US insurance claims data (Truven CCAE)
- 12 month washout
- remove people using both drugs
- remove people with a prior history of stroke
- restricted to people with a diagnosis of major depressive disorder and no prior diagnosis of bipolar disorder or schizophrenia
- Propensity score: regularized logistic regression with treatment as dependent variable, and used 58,285 covariates
- HOI is Stroke: first hospitalization with ICD-9 433,434, or 436 (but then coded as standard concepts)
- fixed-time Cox regression using 10 PS strata





# OHDSI recommendations for evidence generation

- ✓ Post protocol online
  - Prespecify research objectives and design decisions
  
- ✓ Make study code open source
  - From CDM to hazard ratios
  
- ✓ Use validated software
  - OHDSI Methods Library uses unit tests and simulation
  
- ✓ Replicate across several databases
  - 4 included so far, more will follow

<https://github.com/OHDSI/StudyProtocols/LargeScalePopEst>



# Population-level effect estimation



- How do we know the evidence is reliable?



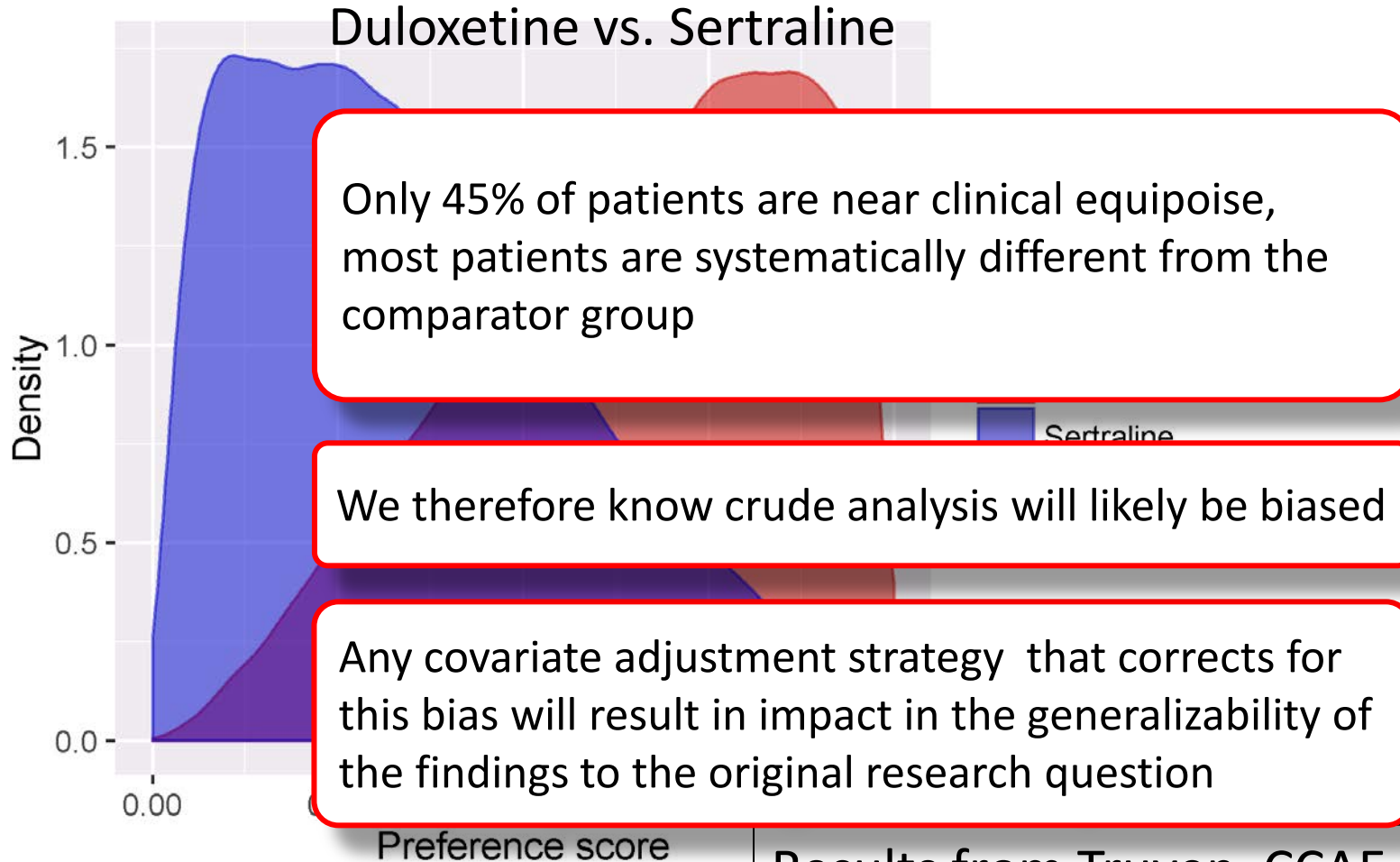
# Standard diagnostics

Most study designs have diagnostics that could be used, e.g.

- Propensity score distribution overlap
- Covariate balance



# Diagnose the propensity score distribution

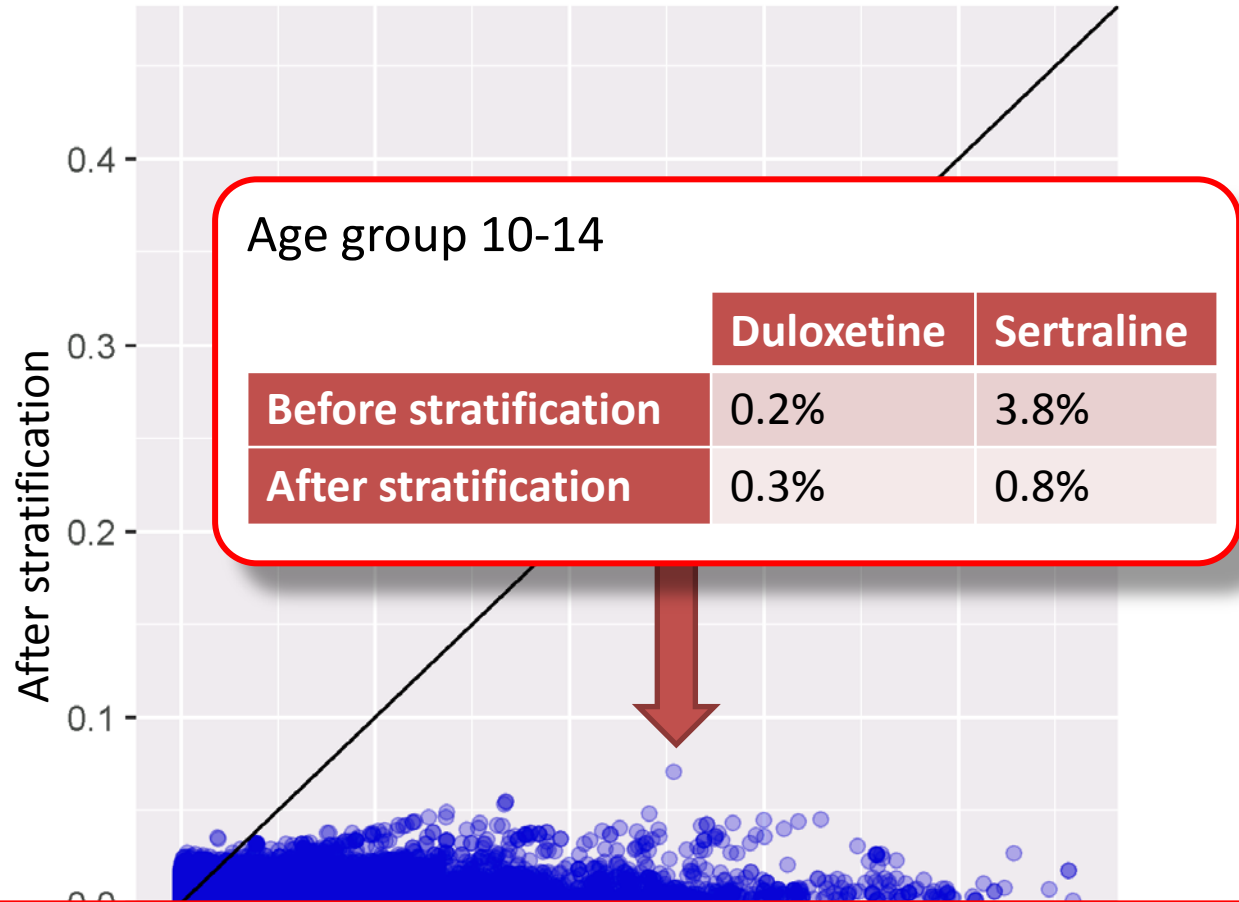


Results from Truven CCAE  
Duloxetine: n = 90,043  
Sertraline: n = 175,950



# Diagnose covariate balance

Standardized difference of mean



After stratification on the propensity score, all 58,285 covariates have standardized difference of mean  $< 0.1$



# Empirical evaluation of the study

- Control  
exposure-outcome for which the effect size is known
- Negative control  
exposure-outcome where relative risk is believed to be 1
- Negative controls for comparative effectiveness  
outcomes not believed to be caused by either treatments

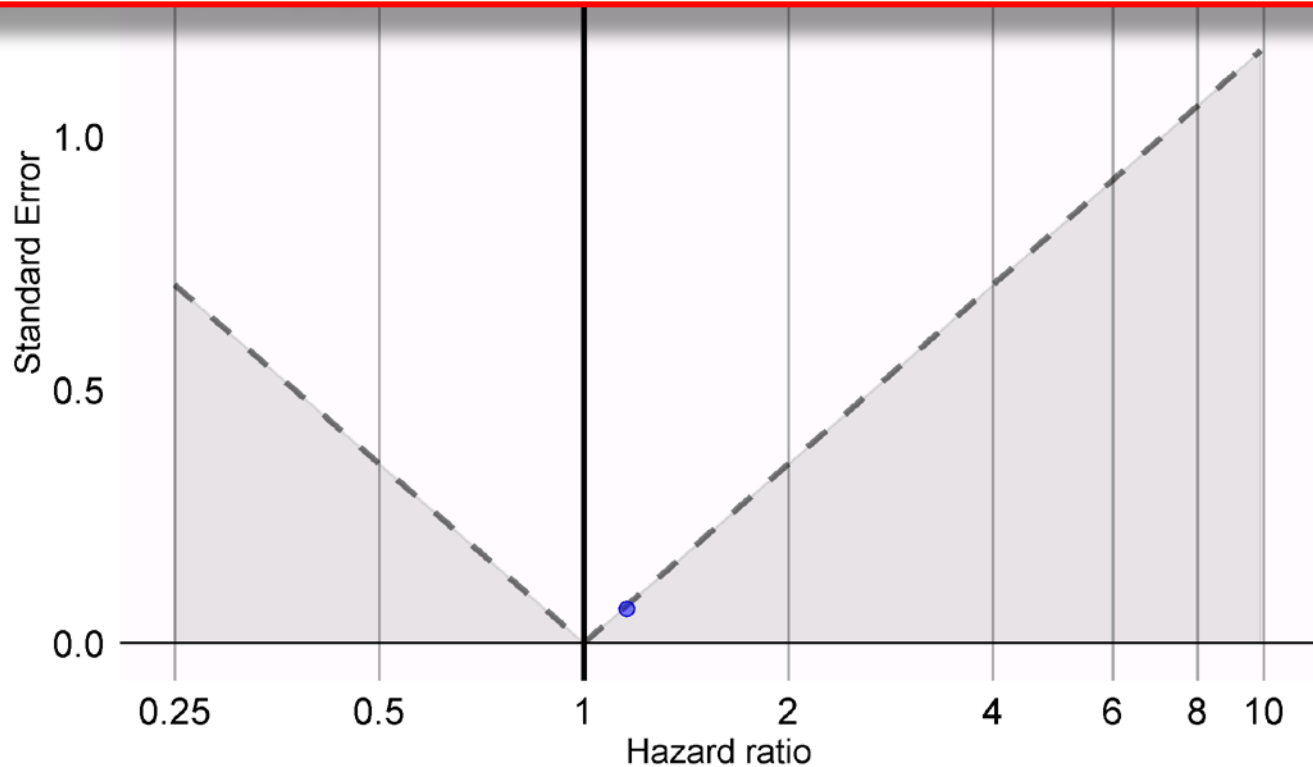
Example: ingrowing nail



# Negative control: ingrowing nail

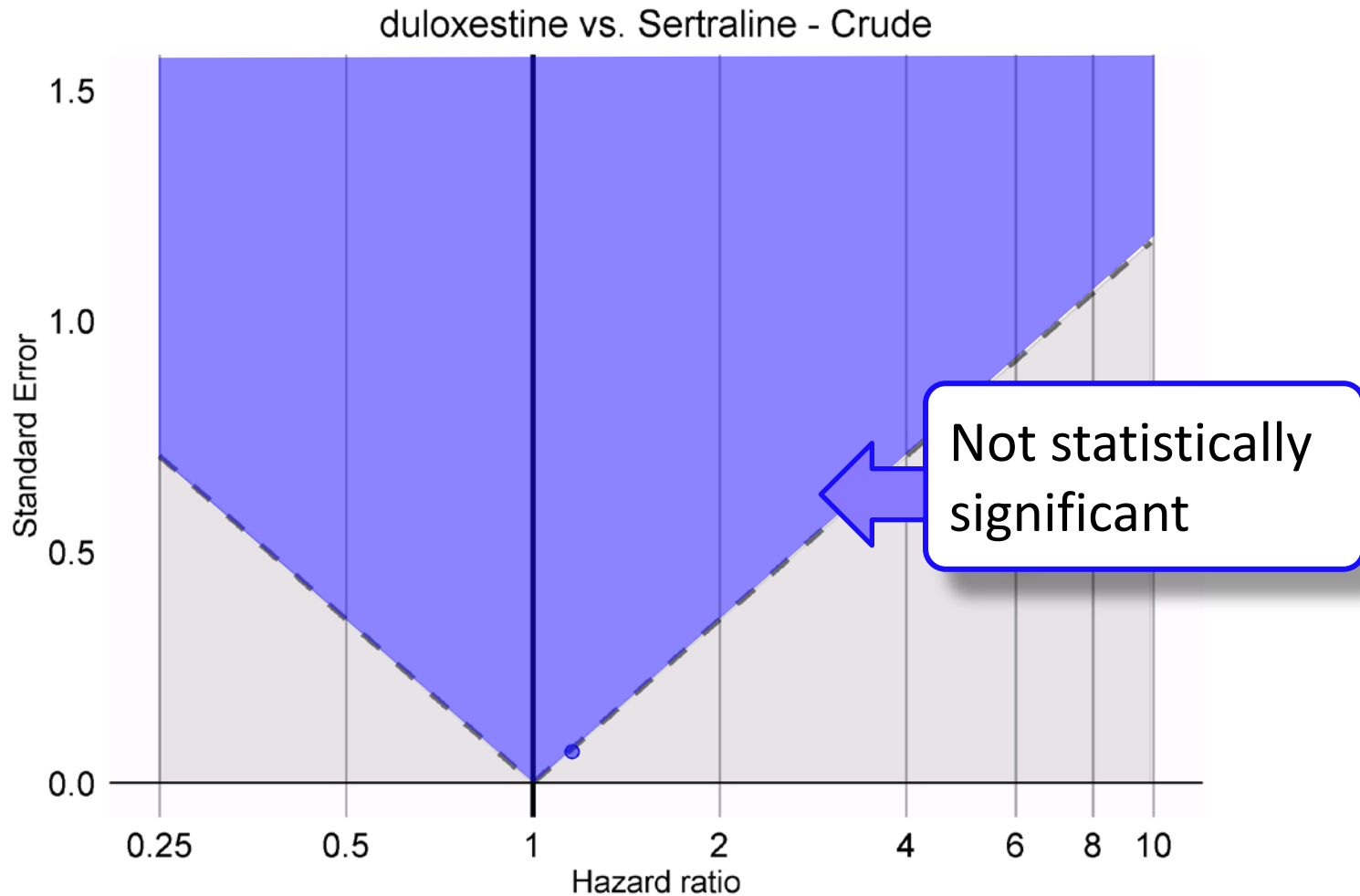
Crude estimate:

HR = 1.16 (1.01 – 1.32),  $p = 0.03$





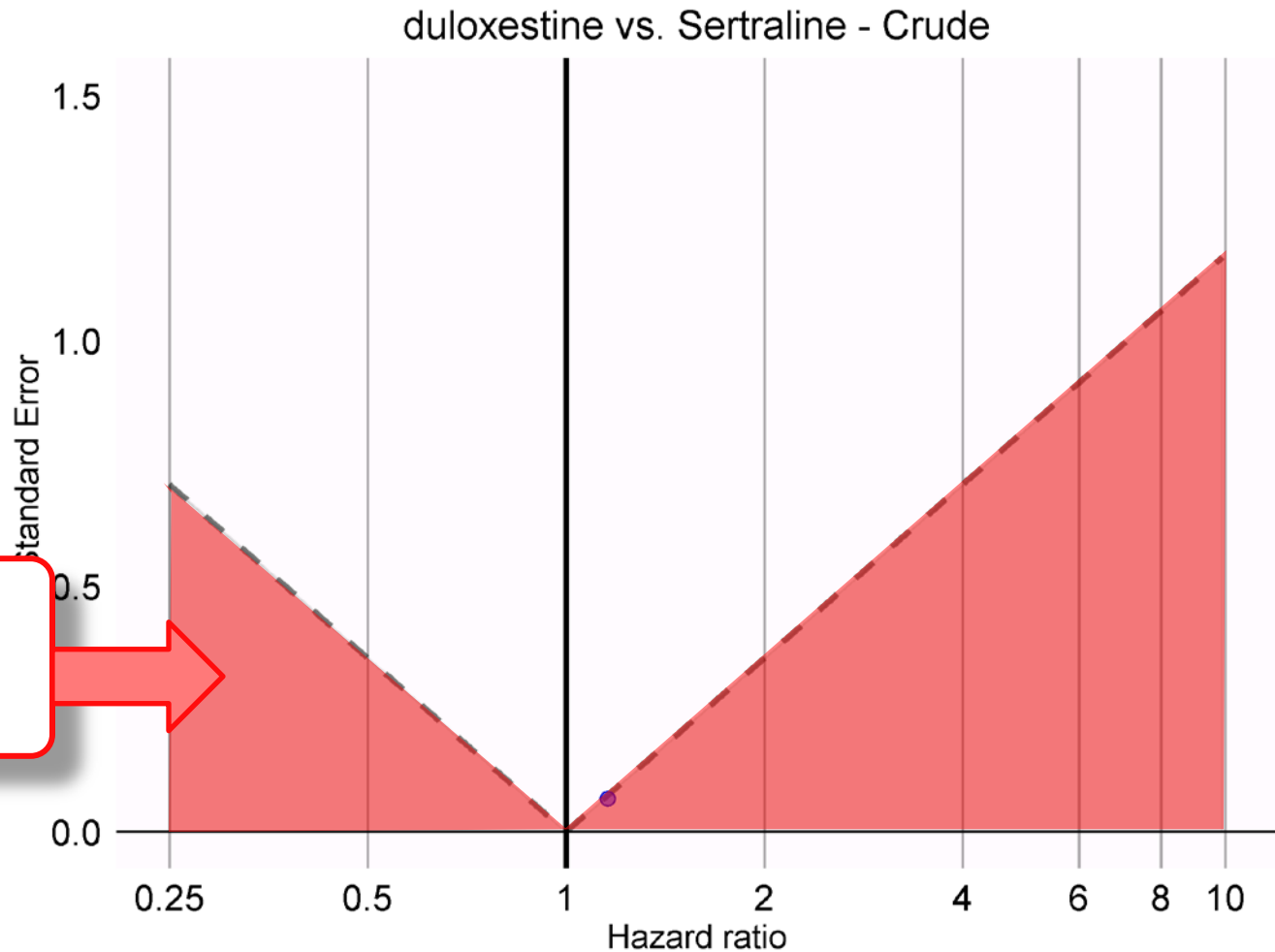
# Negative control: ingrowing nail







# Negative control: ingrowing nail



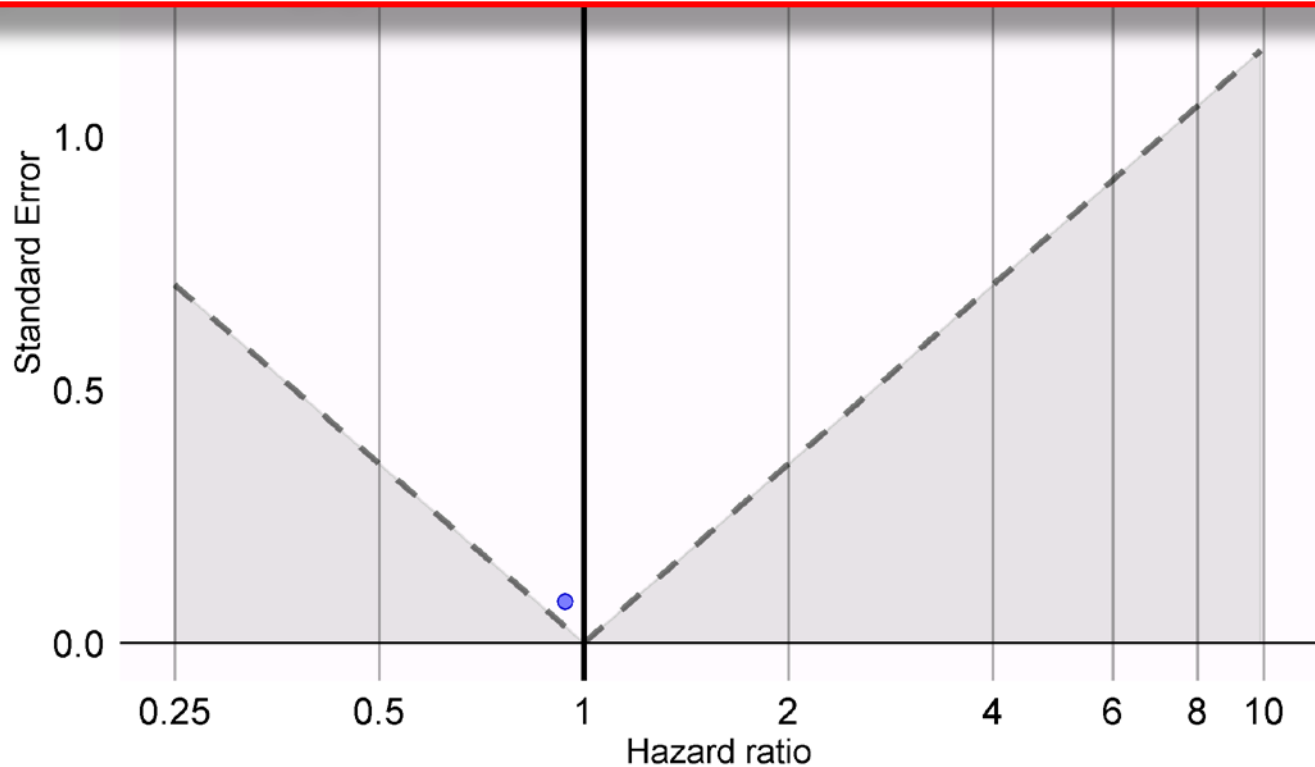
Statistically significant



# Negative control: ingrowing nail

Adjusted estimate:

HR = 0.94 (0.80 – 1.10),  $p = 0.44$





# Depression – negative controls

Acariasis	Ingrowing nail
Amyloidosis	Iridocyclitis
Ankylosing spondylitis	Irritable bowel syndrome
Aseptic necrosis of bone	Lesion of cervix
Astigmatism	Lyme disease
Bell's palsy	Malignant neoplasm of endocrine gland
Benign epithelial neoplasm of skin	Mononeuropathy
Chalazion	Onychomycosis
Chondromalacia	Osteochondropathy
Crohn's disease	Paraplegia
Croup	Polyp of intestine
Diabetic oculopathy	Presbyopia
Endocarditis	Pulmonary tuberculosis
Endometrial hyperplasia	Rectal mass
Enthesopathy	Sarcoidosis
Epicondylitis	Scar
Epstein-Barr virus disease	Seborrheic keratosis

Generated with the help of LAERTES (see posters)

Hemangioma	Periostitis
Hodgkin's disease	Toxic goiter
Human papilloma virus infection	Ulcerative colitis
Hypoglycemic coma	Viral conjunctivitis
Hypopituitarism	Viral hepatitis
Impetigo	Visceroptosis

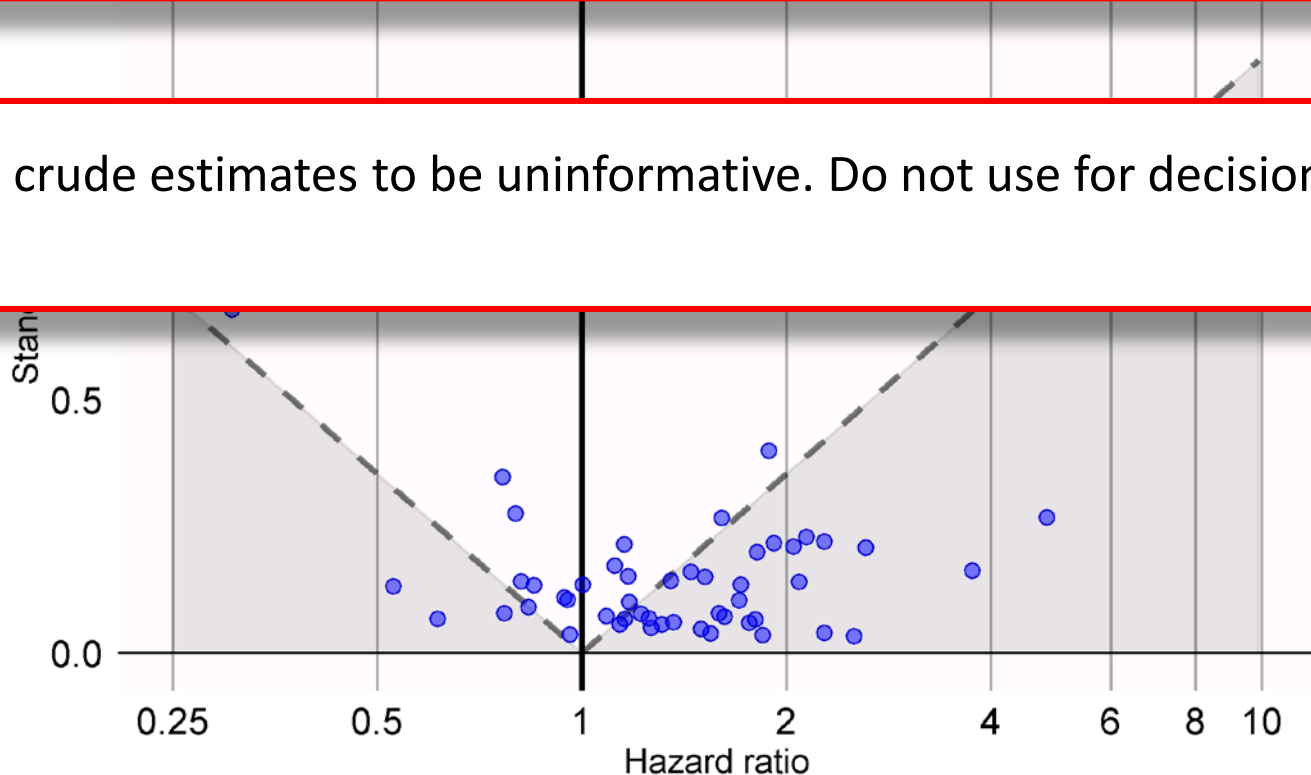


# All negative controls - crude

We would expect 5% of negative controls to have  $p < 0.05$

Instead, 68% has  $p < 0.05$ !

We found crude estimates to be uninformative. Do not use for decision making!



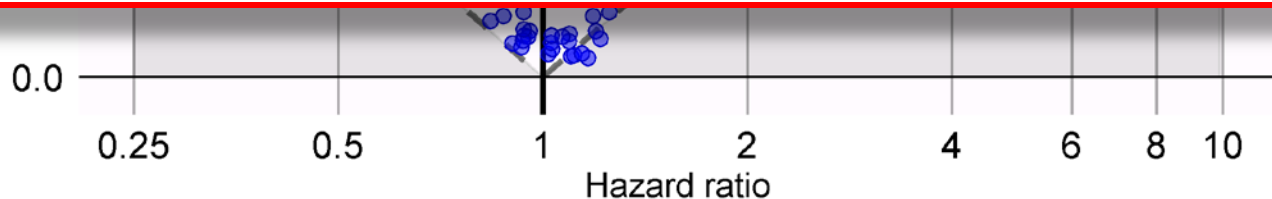


# All negative controls - adjusted

When using the propensity score, 16% have  $p < 0.05$

In the past, we've shown you how you can perform p-value calibration:

- P-value represents probability of estimate when true RR = 1
- Negative controls provide empirical distribution of estimates when RR = 1
- Use empirical null distribution to compute calibrated p-value





# P-value calibration

duloxetine vs. Sertraline - Adjusted

1.5

After calibration, 4% have  $p < 0.05$  (was 16%)

What if  $HR < 1$ ?

Standard

0.5

0.0

Calibrated  $p < 0.05$

0.25

0.5

1

2

4

6

8

10

Hazard ratio



# Trouble with positive controls

- Often very few positive examples for a particular comparison
- Exact effect size never known with certainty (and depends on population)
- Doctors also know they're positive, and will change behavior accordingly

Drug Saf (2014) 37:655–659  
DOI 10.1007/s40264-014-0198-z

CURRENT OPINION

## **Zoo or Savannah? Choice of Training Ground for Evidence-Based Pharmacovigilance**

G. Niklas Norén · Ola Caster ·  
Kristina Juhlin · Marie Lindquist



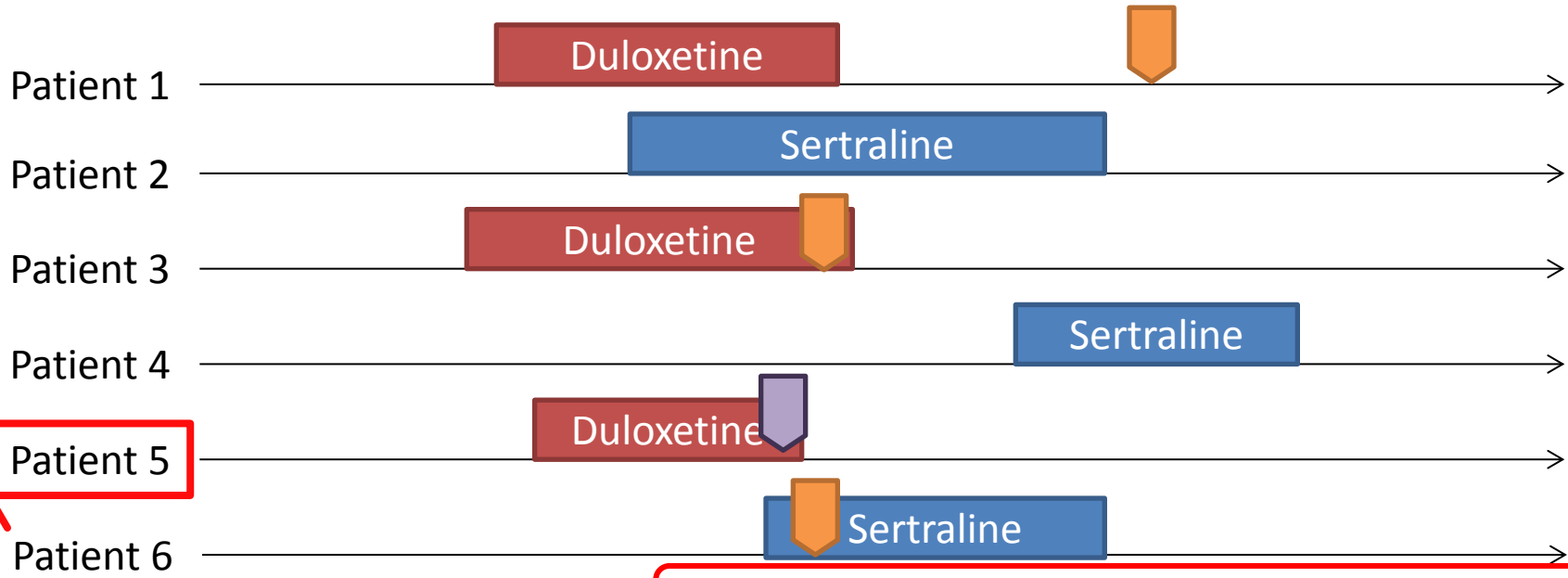
# Creating positive controls

- Start with negative controls:  $RR = 1$
- Add simulated outcomes during exposure until desired  $RR$  is achieved
- Injected outcomes should behave like ‘real’ outcomes: preserve confounding structure by injecting outcomes for people at high risk







# Creating positive controls



Patient 5

New RR = 2 (but with same confounding)

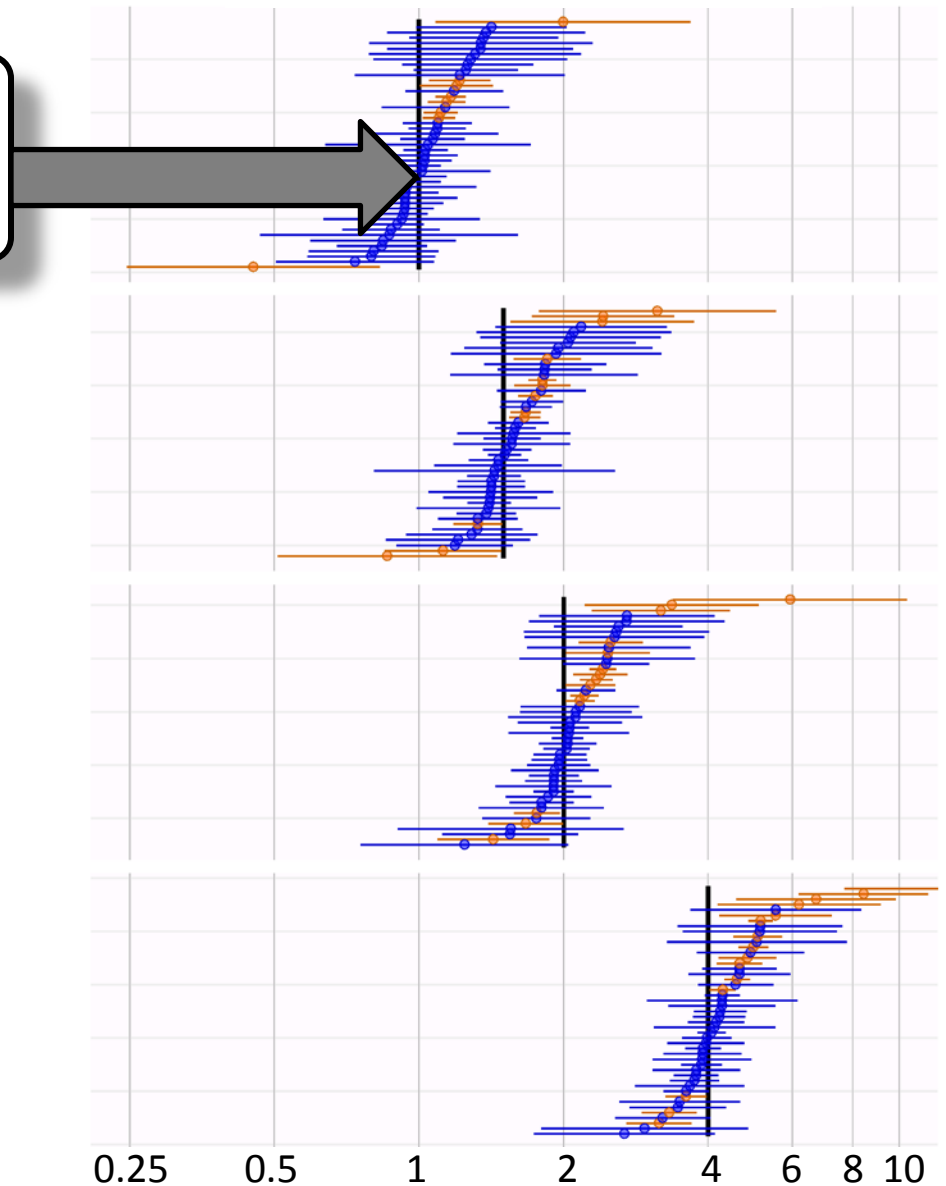
-  Ingrowing nail
-  Injected ingrowing nail

Predictive model of outcome indicates this is a high-risk patient



# Estimated effects for positive controls

Black line indicates true hazard ratio





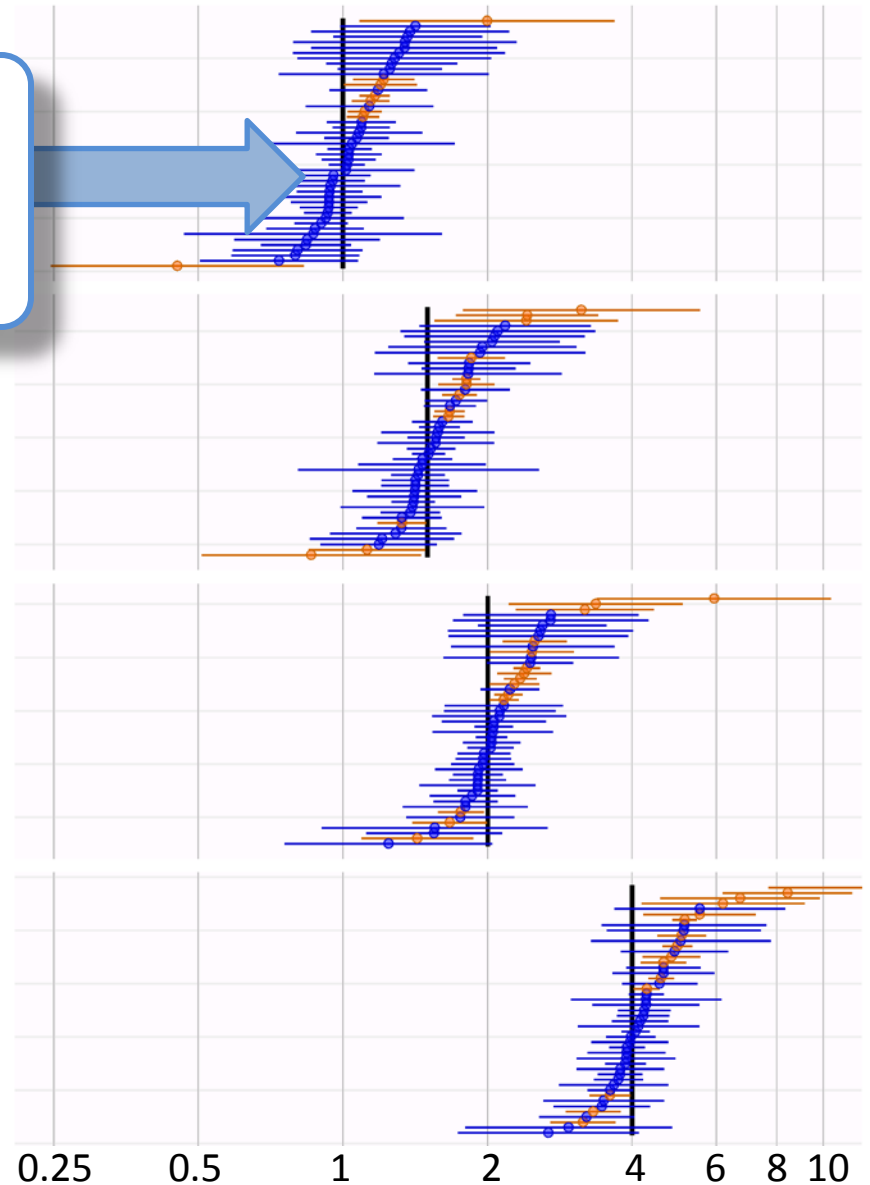
# Estimating effects for positive controls

Ingrowing nail

True RR = 1

Estimated RR = 0.94 (0.80 – 1.10)

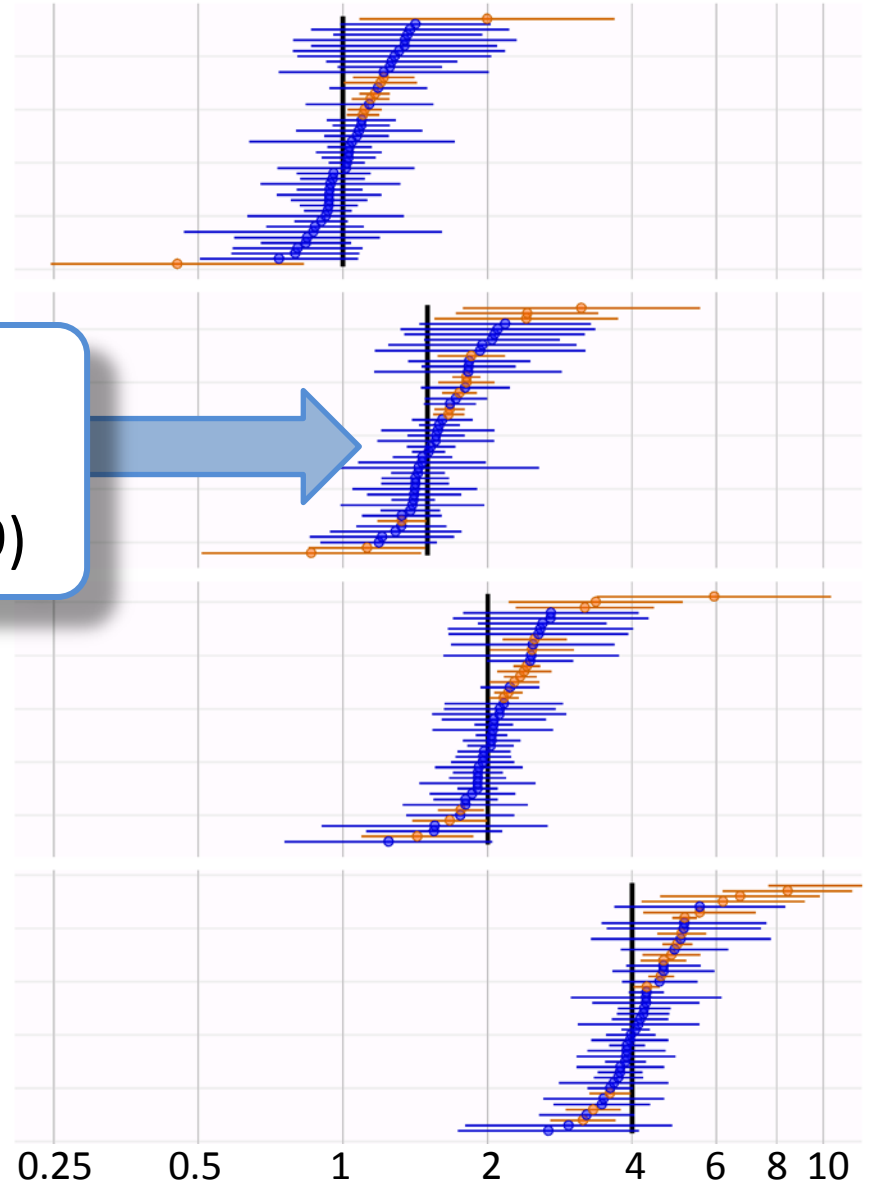
duloxetine vs. Sertraline - Adjusted





# Estimating effects for positive controls

duloxetine vs. Sertraline - Adjusted



Ingrowing nail+

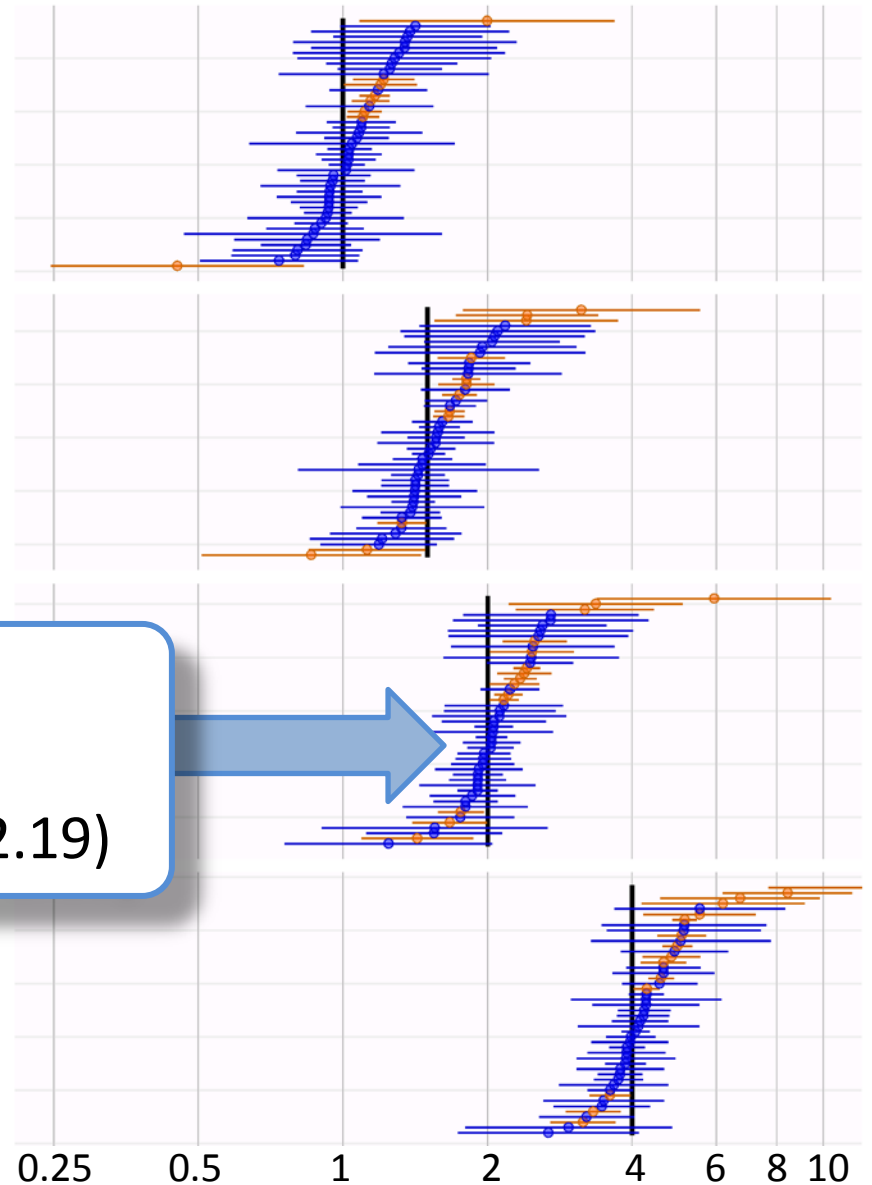
True RR = 1.5

Estimated RR = 1.47 (1.27 – 1.69)



# Estimating effects for positive controls

duloxetine vs. Sertraline - Adjusted



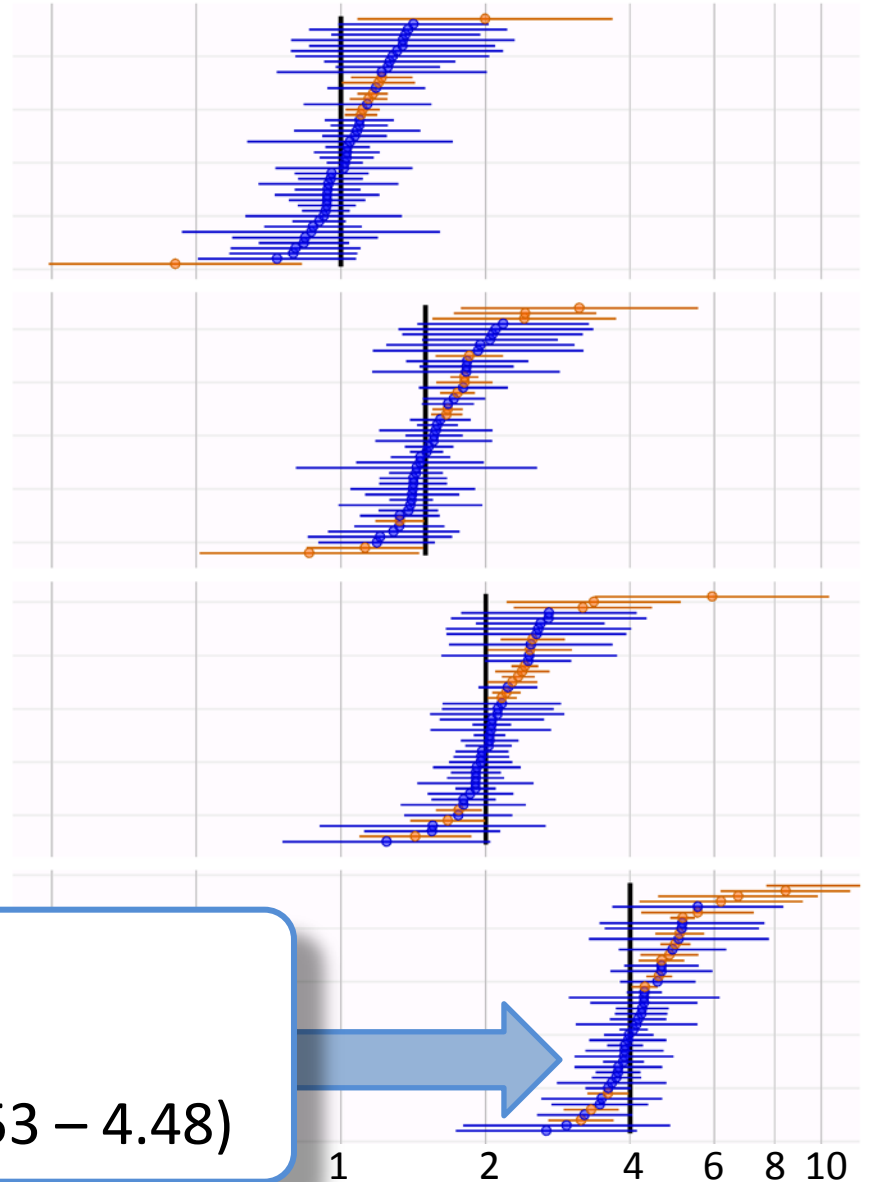
Ingrowing nail++

True RR = 2

Estimated RR = 1.91 (1.67 – 2.19)

# Estimating effects for positive controls

duloxetine vs. Sertraline - Adjusted



Ingrowing nail+++

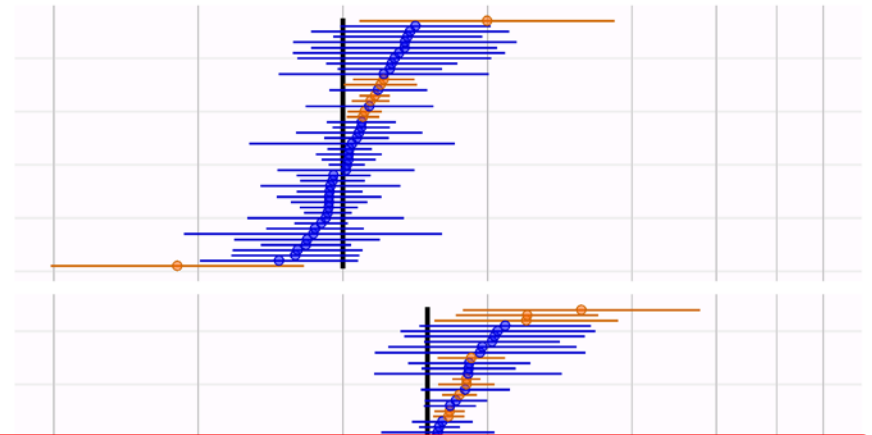
True RR = 4

Estimated RR = 3.89 (3.53 – 4.48)

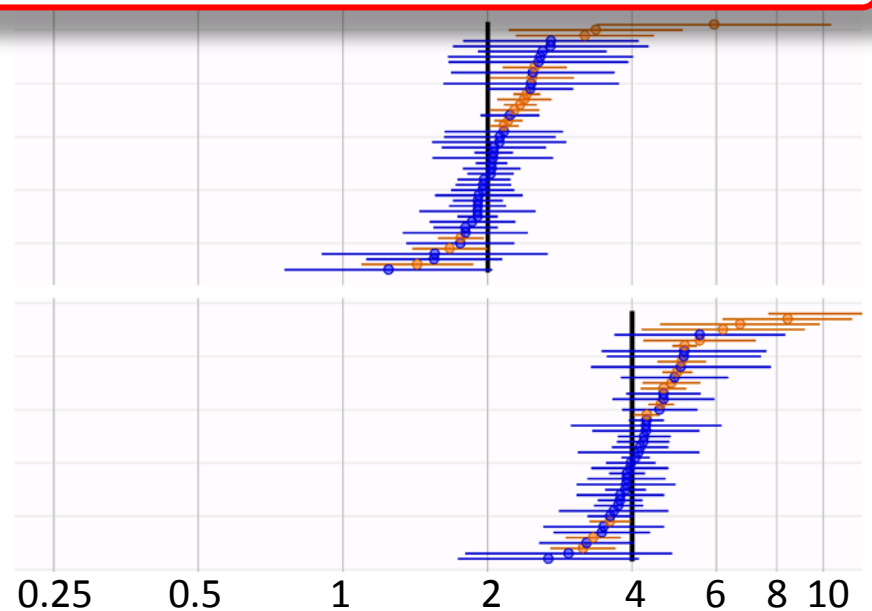


# Estimating effects for positive controls

duloxetine vs. Sertraline - Adjusted



Analysis suggests bias remains constant with effect size





# Evaluating coverage of the CI

Coverage

duloxetine vs. Sertraline - Adjusted

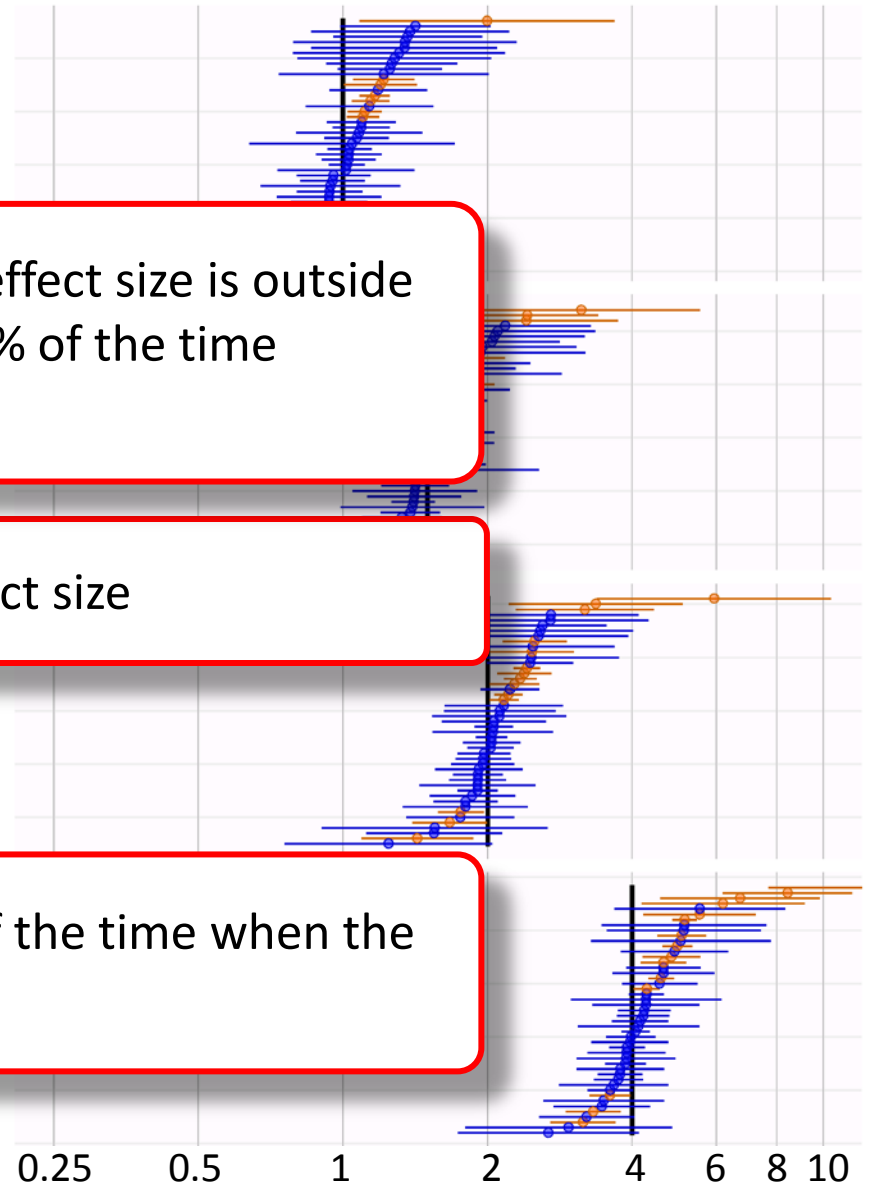
83%

Coverage of 83% means the true effect size is outside of the 95% confidence interval 17% of the time (when the true RR = 1)

Coverage decreases with true effect size

70%

Missing the true effect size 30% of the time when the true RR = 2!

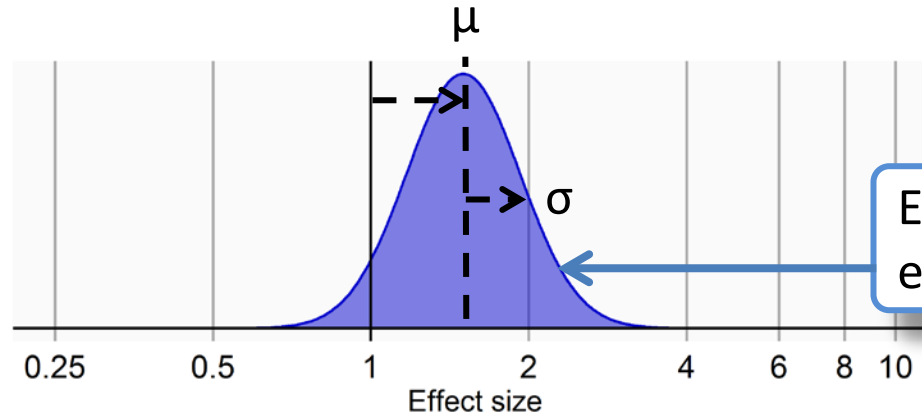




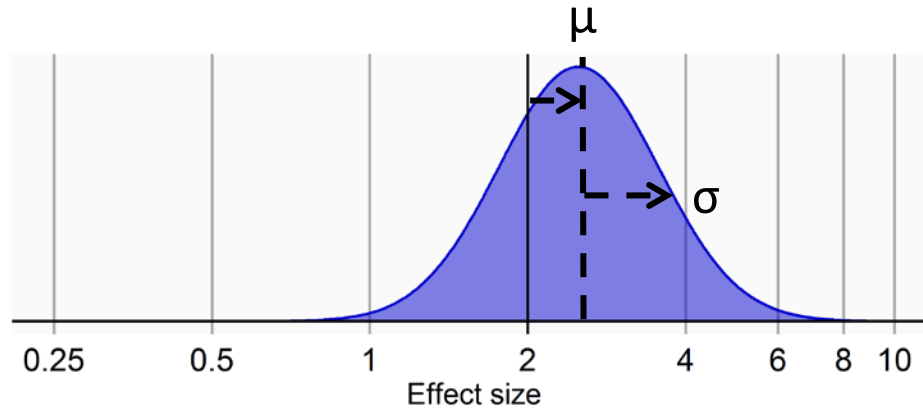


# Confidence interval calibration

$HR_{true} = 1$



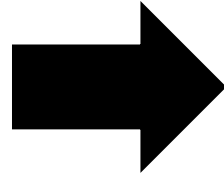
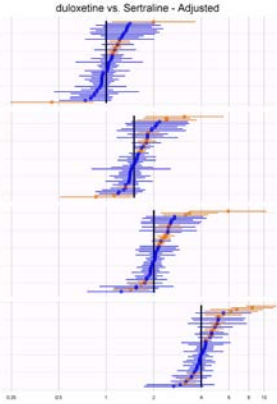
$HR_{true} = 2$



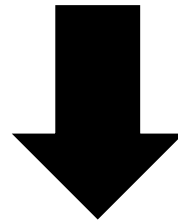
$$\mu = \alpha_{\mu} + \beta_{\mu} \log(HR_{true})$$

$$\sigma = \alpha_{\sigma} + \beta_{\sigma} \log(HR_{true})$$

# Calibrating a confidence interval



$$\mu = 0.04 + 1.01 \log(\text{HR}_{\text{true}})$$
$$\sigma = 0.07 + 0.05 \log(\text{HR}_{\text{true}})$$

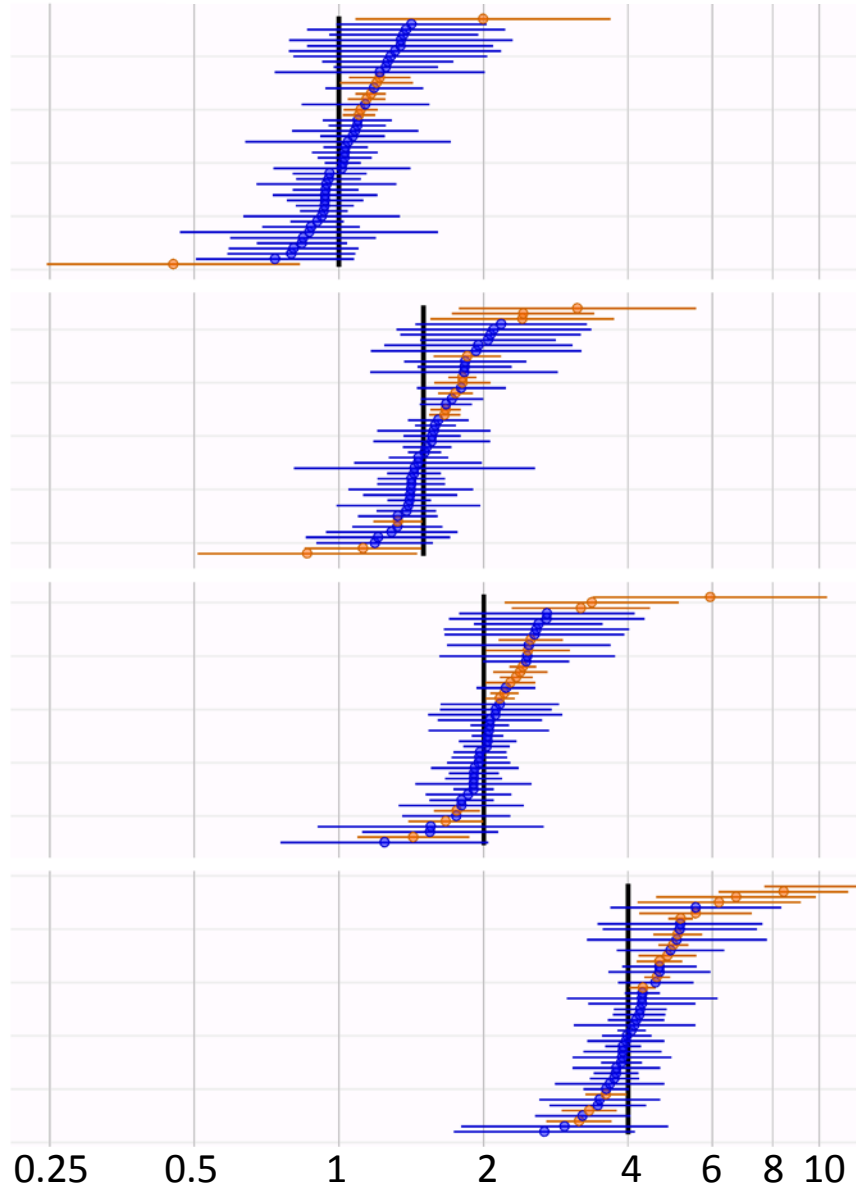


Confidence intervals were too narrow, so made wider to get to nominal coverage

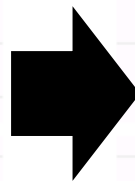
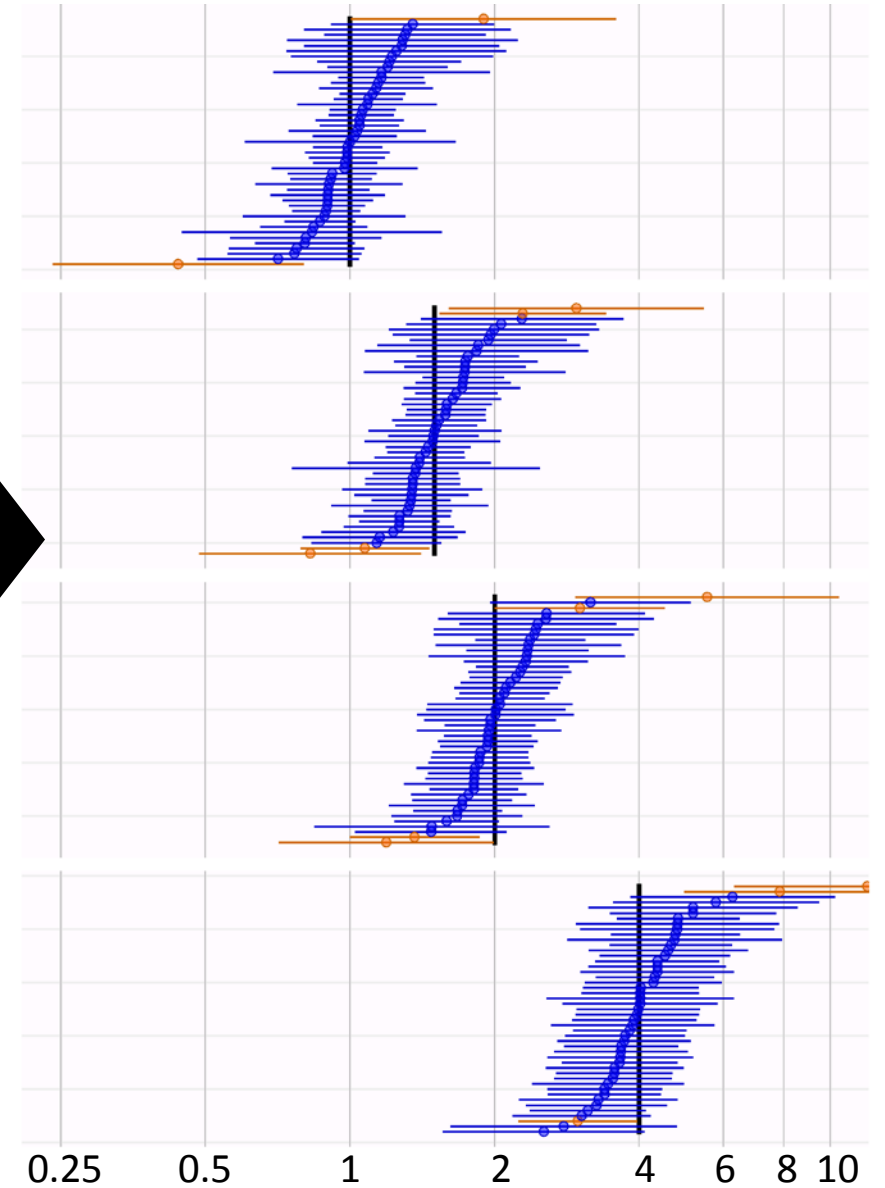


# Confidence interval calibration

Uncalibrated



Calibrated





# Confidence interval calibration

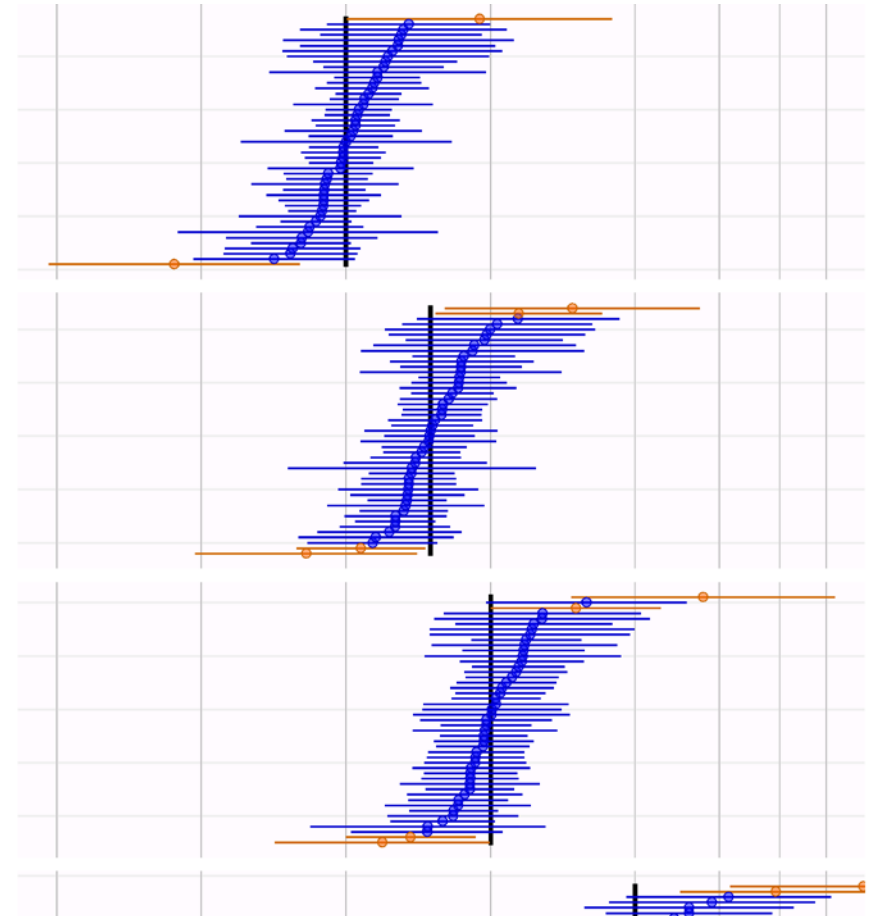
Coverage

Calibrated

96%

91%

91%



Confidence interval calibration complements p-value calibration

0.25 0.5 1 2 4 6 8 10



# Current evidence for stroke

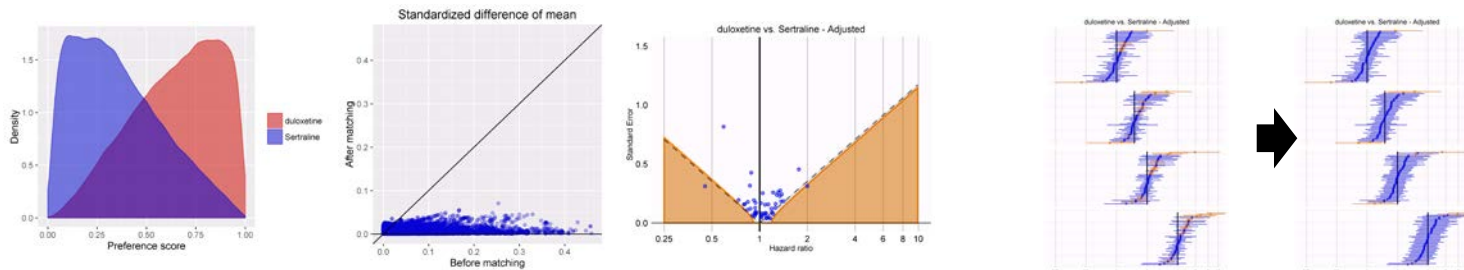
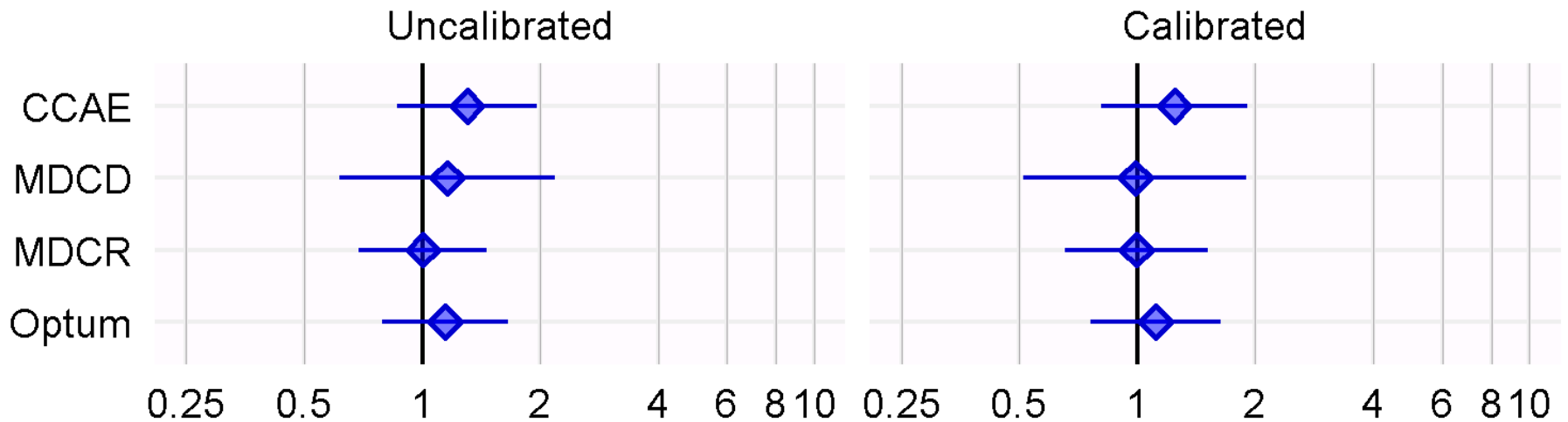
Result from Lee et al.

	Crude Hazard Ratio (95% CI)	<i>P</i>	Adjusted Hazard Ratio <sup>a</sup> (95% CI)	<i>P</i>
<b>Main analyses</b>				
SNRIs (n= 76,920) vs SSRIs (n= 582,650)				
Ischemic stroke	0.92 (0.83–1.02)	.12	1.01 (0.90–1.12)	.91



# Proposed evidence for stroke

## Duloxetine vs. Sertraline



Results are comparable to Lee et al., but we provide the context to interpret the results

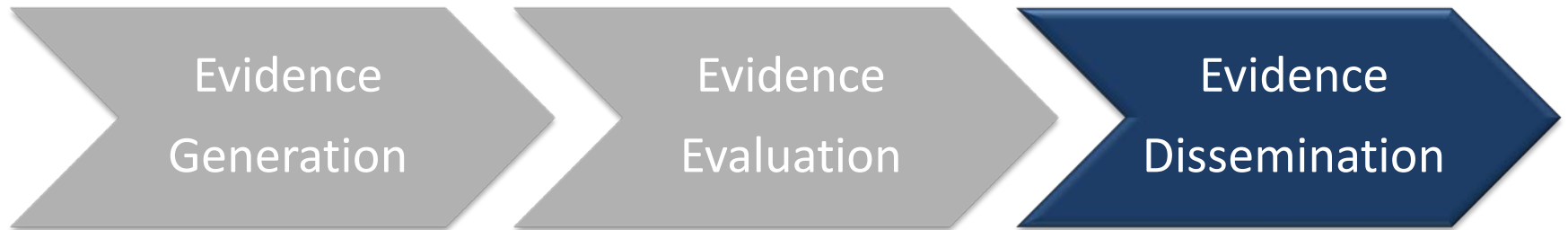


# OHDSI recommendations for evidence evaluation

- ✓ Produce standard diagnostics
  - E.g. for cohort studies diagnose the propensity score distribution, covariate balance, etc.
  
- ✓ Include negative controls
  - Estimate the error when the null is true
  
- ✓ Create positive controls
  - Estimate the error when  $RR > 1$
  
- ✓ Calibrate p-value and confidence intervals
  - Restoring nominal characteristics



# Population-level effect estimation



- How do we share evidence to inform decision making?





# Evidence dissemination

- Traditionally, this evidence is disseminated through the scientific literature
- How well does that work?

# Automated extraction of effect sizes from literature



**RESULTS:** In comparison with distant past users of BP, current users of BP showed an almost twofold increased risk of AF: odds ratio (OR) = 1.78 and 95% CI = 1.46-2.16. Specifically, alendronate users were mostly associated with AF as compared with distant past use of BP (OR, 1.97; 95% CI 1.59-2.43).

## Abstract

Bisphosphonate treatment is used to prevent bone fractures. A controversial association of bisphosphonate use and risk of atrial fibrillation has been reported. In our study, current alendronate users were associated with a higher risk of atrial fibrillation as compared with those who had stopped bisphosphonate (BP) therapy for more than 1 year.

**INTRODUCTION:** Bisphosphonates are widely used to prevent bone fractures. Controversial findings regarding the association between bisphosphonate use and the risk of atrial fibrillation (AF) have been reported. The aim of this study was to evaluate the risk of AF in association with BP exposure.

**METHODS:** We performed a nested case-control study using the databases of drug-dispensing and hospital discharge diagnoses from five Italian regions. The data cover a period ranging from July 1, 2003 to December 31, 2006. The study population comprised new users of bisphosphonates aged 55 years and older. Patients were followed from the first BP prescription until an occurrence of an AF diagnosis (index date, i.e., ID), cancer, death, or the end of the study period, whichever came first. For the risk estimation, any AF case was matched by age and sex to up to 10 controls from the same source population. A conditional logistic regression was performed to obtain the odds ratio with 95% confidence intervals (CI). The BP exposure was classified into current (<90 days prior to ID), recent (91-180), past (181-364), and distant past ( $\geq 365$ ) use, with the latter category being used as a reference point. A subgroup analysis by individual BP was then carried out.

**RESULTS:** In comparison with distant past users of BP, current users of BP showed an almost twofold increased risk of AF: odds ratio (OR) = 1.78 and 95% CI = 1.46-2.16. Specifically, alendronate users were mostly associated with AF as compared with distant past use of BP (OR, 1.97, 95% CI, 1.59-2.43).

**CONCLUSION:** In our nested case-control study, current users of BP are associated with a higher risk of atrial fibrillation as compared with those who had stopped BP treatment for more than 1 year.

PMID: 25752621 [PubMed - indexed for MEDLINE] PMCID: PMC4428862 [Free PMC Article](#)



Images from this publication. [See all images \(1\)](#) [Free text](#)

## Similar articles

Oral bisphosphonates and risk of ischemic stroke: a case-control study [Osteoporos Int. 2011]

Assessing the risk of osteonecrosis of the jaw due to bisphosphonate therapy [Osteoporos Int. 2013]

Use of bisphosphonate and risk of atrial fibrillation in older women [Osteoporos Int. 2012]

**Review** Bisphosphonates and atrial fibrillation: systematic review and meta-analysis [Drug Saf. 2009]

**Review** Risk of atrial fibrillation with use of oral and intravenous bisphosphonates [Am J Cardiol. 2014]

[See reviews...](#)

[See all...](#)

## Related information

Articles frequently viewed together

MedGen

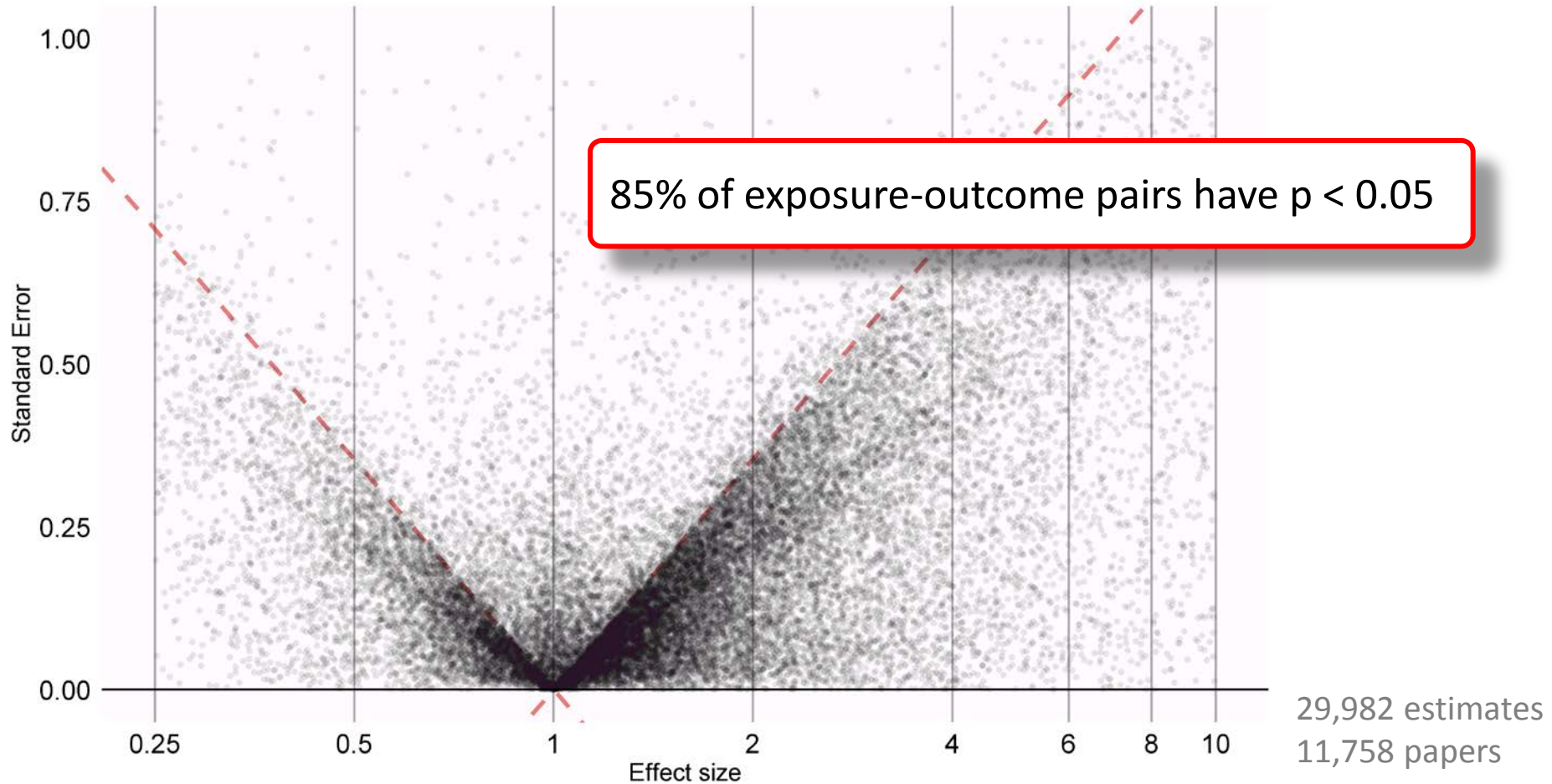
References for this PMC Article

Free in PMC

## Recent Activity



# Observational research results in literature





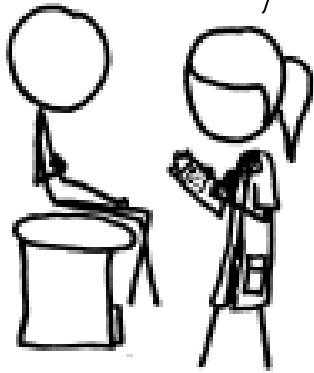
# What went wrong?

- Observational study bias
- Publication bias
- P-hacking

# Observational study bias

I have a headache and my stomach really hurts!

I'll prescribe drug A for your headache, it's safe for people at risk of stomach bleeding.



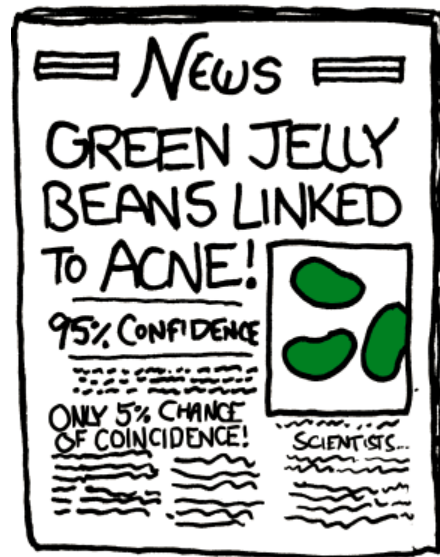
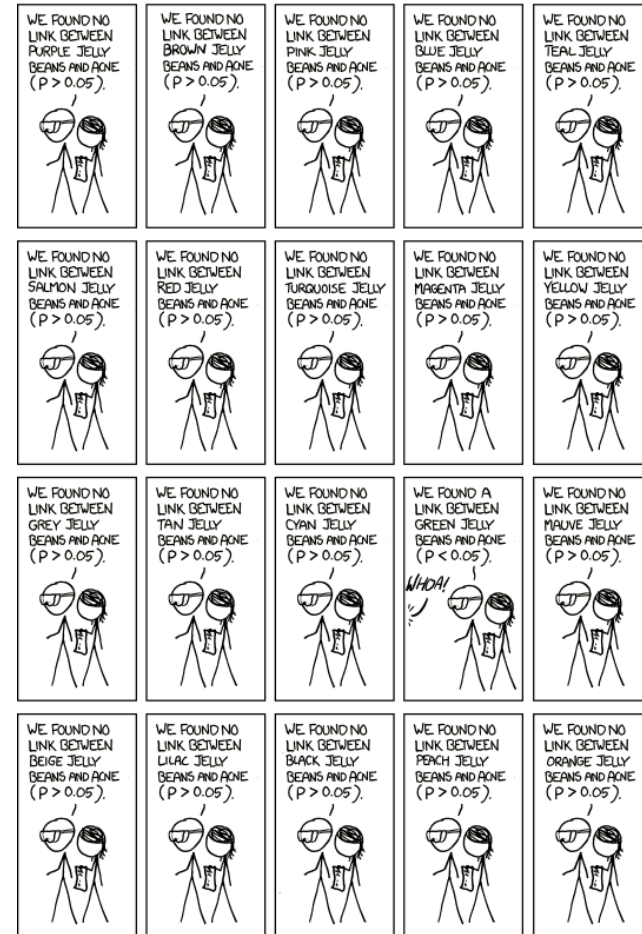
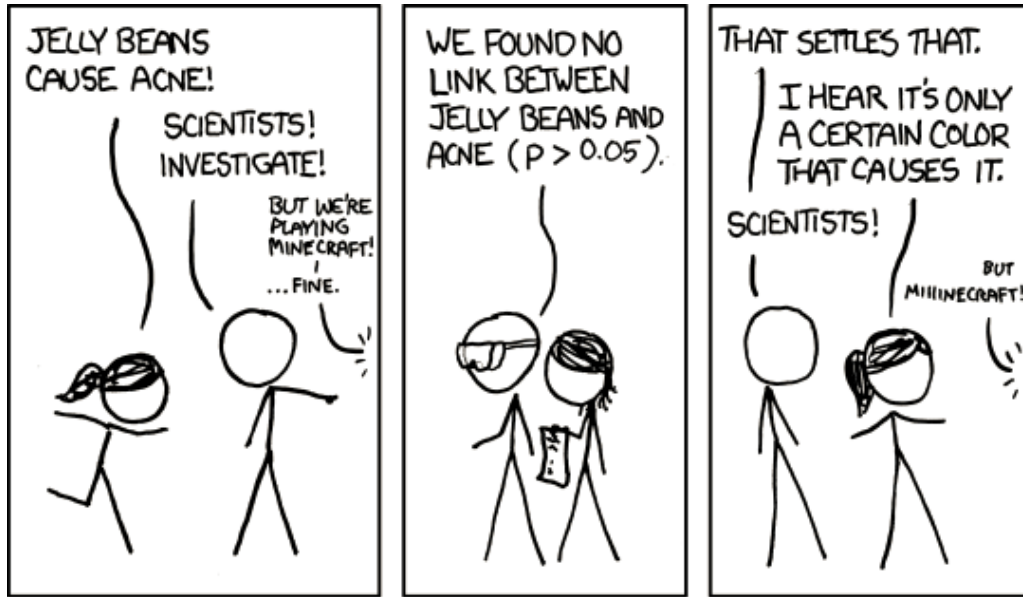
One week later...

I took drug A, now I have a stomach bleeding!

Ha! Drug A causes stomach bleedings!



# Publication bias



# P-hacking

PhD Student!

I think A may cause B,  
go investigate!

Yes professor!



I ran the analysis:  
 $p > .05$

But did you adjust  
for confounder Z?

Ehh, no

Let me get  
right back to you



After adjustment  
for Z,  $p < .05!$

Yay! Lets publish  
a paper!





# A solution?

Stop doing one study at a time!





# What if we considered all outcomes?

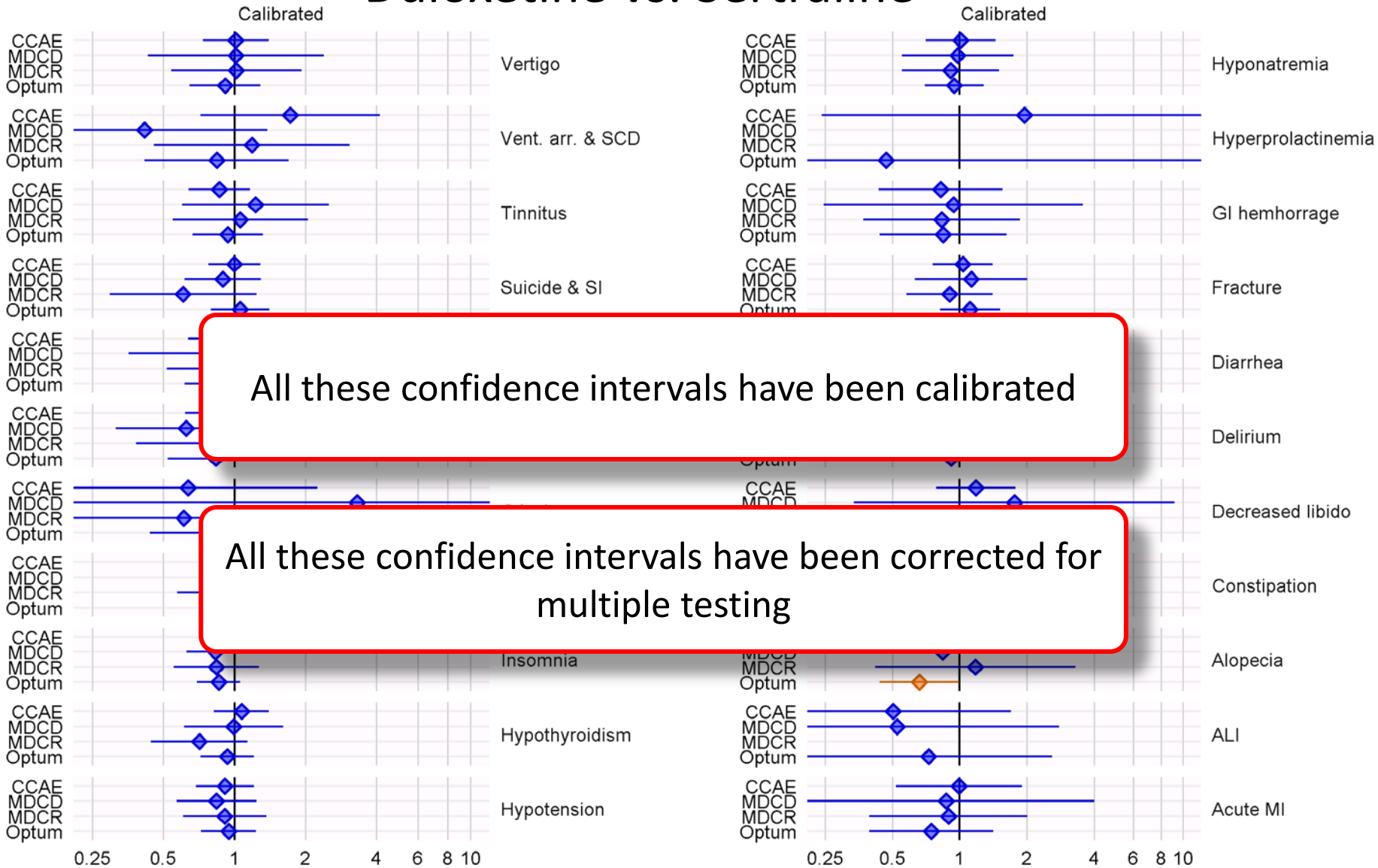
Duloxetine vs. Sertraline for these 22 outcomes:

Acute liver injury	Hypotension
Acute myocardial infarction	Hypothyroidism
Alopecia	Insomnia
Constipation	Nausea
Decreased libido	Open-angle glaucoma
Delirium	Seizure
Diarrhea	Stroke
Fracture	Suicide and suicidal ideation
Gastrointestinal hemorrhage	Tinnitus
Hyperprolactinemia	Ventricular arrhythmia and sudden cardiac death
Hyponatremia	Vertigo



# All outcomes

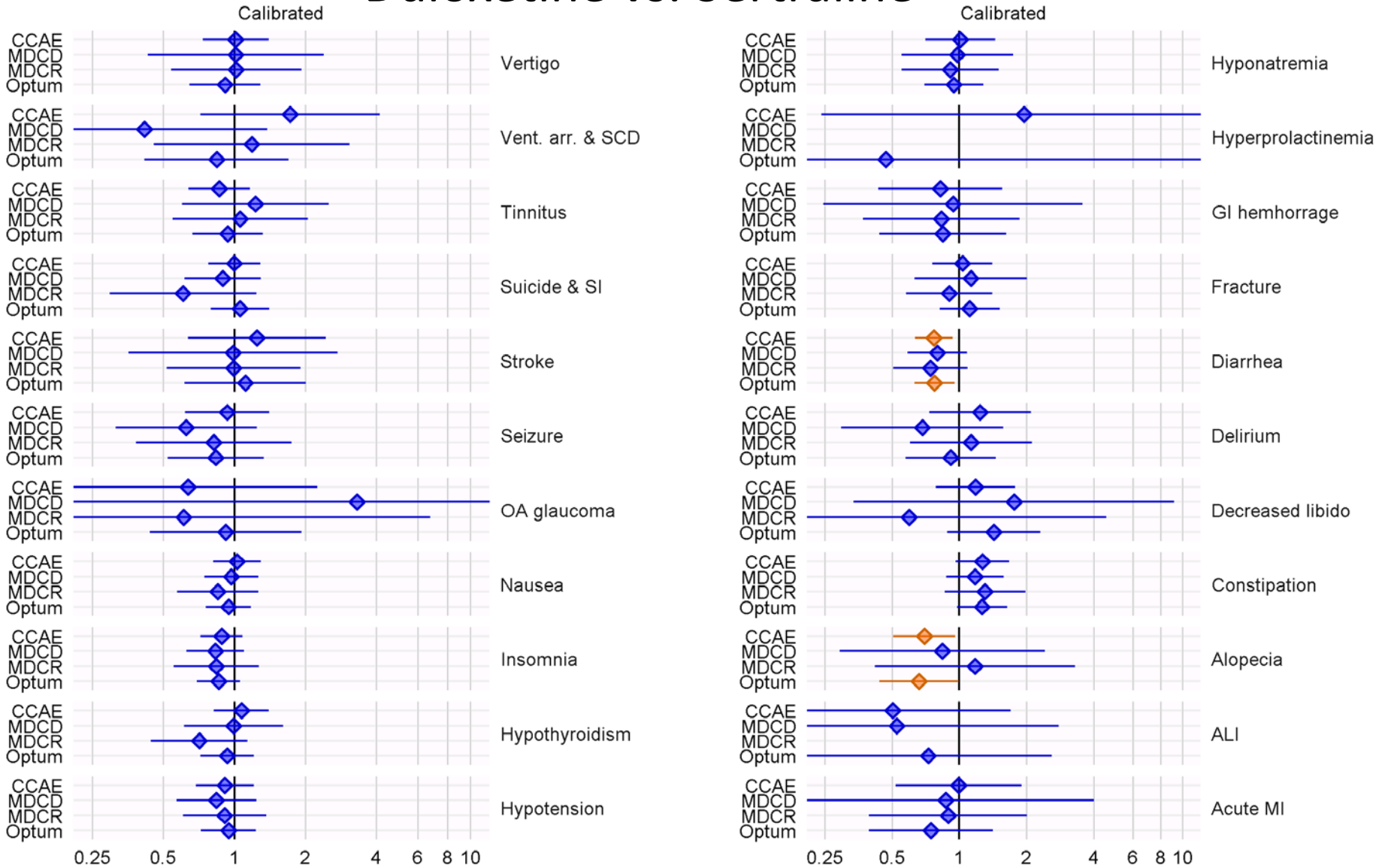
## Duloxetine vs. Sertraline





# All outcomes

## Duloxetine vs. Sertraline





# What if we consider all treatments?

Type	Class	Treatment
Drug	Atypical	Bupropion
Drug	Atypical	Mirtazapine
Procedure	ECT	Electroconvulsive therapy
Procedure	Psychotherapy	Psychotherapy
Drug	SARI	Trazodone
Drug	SNRI	Desvenlafaxine
Drug	SNRI	duloxetine
Drug	SNRI	venlafaxine
Drug	SSRI	Citalopram
Drug	SSRI	Escitalopram
Drug	SSRI	Fluoxetine
Drug	SSRI	Paroxetine
Drug	SSRI	Sertraline
Drug	SSRI	vilazodone
Drug	TCA	Amitriptyline
Drug	TCA	Doxepin
Drug	TCA	Nortriptyline



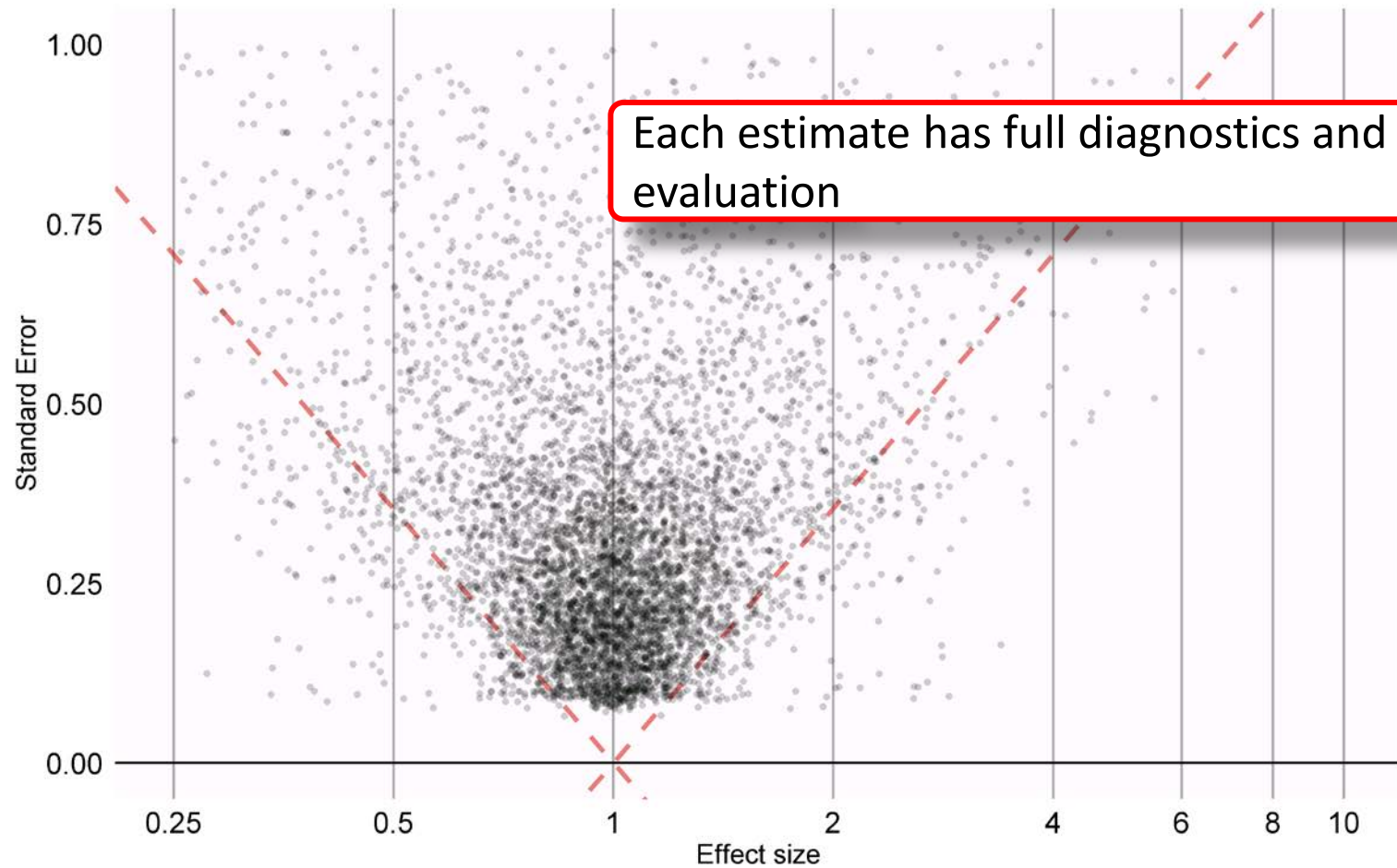
# Large-scale estimation for depression

- **17 treatments**
- $17 * 16 = 272$  comparisons
- **22 outcomes**
- $272 * 22 = 5,984$  effect size estimates
- **4 databases** (Truven CCAE, Truven MDCCD, Truven MDCCR, Optum)
- $4 * 5,984 = \mathbf{23,936}$  estimates



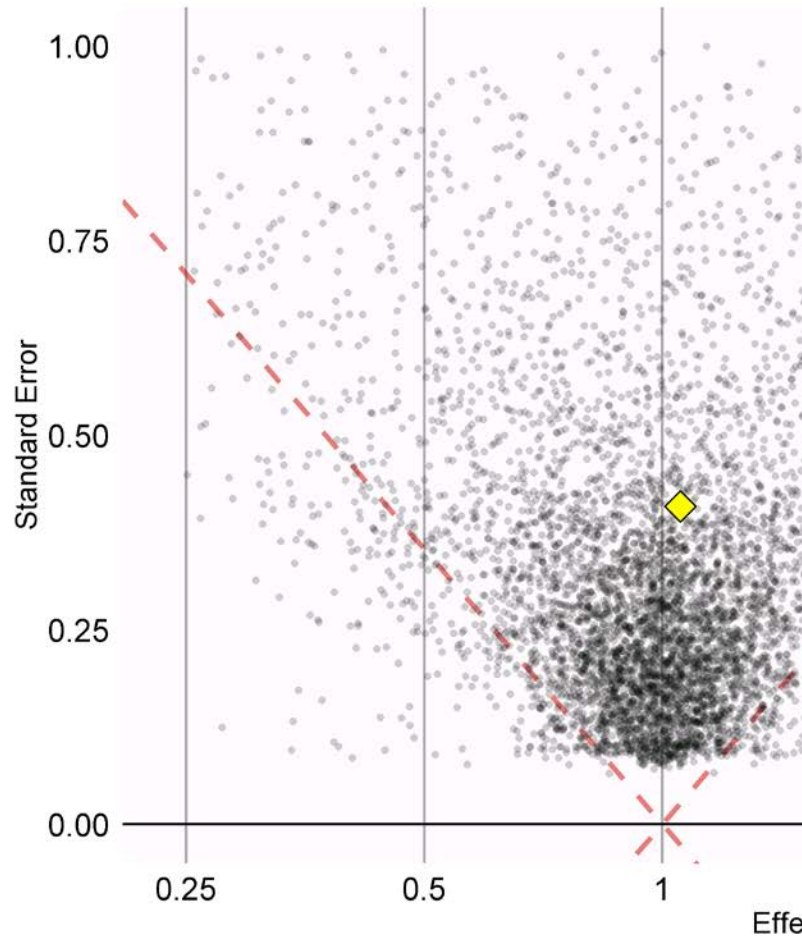


# Large-scale estimation for depression



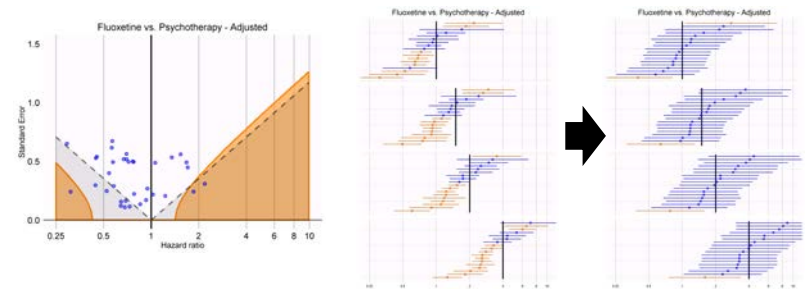
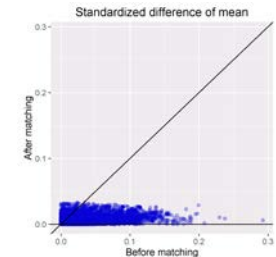
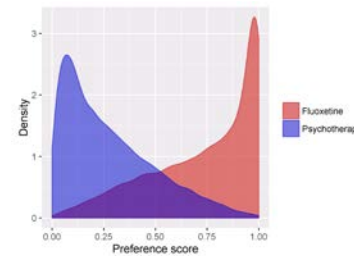


# Example 1



Fluoxetine vs. psychotherapy  
Suicide ideation  
Database: Truven MDCR

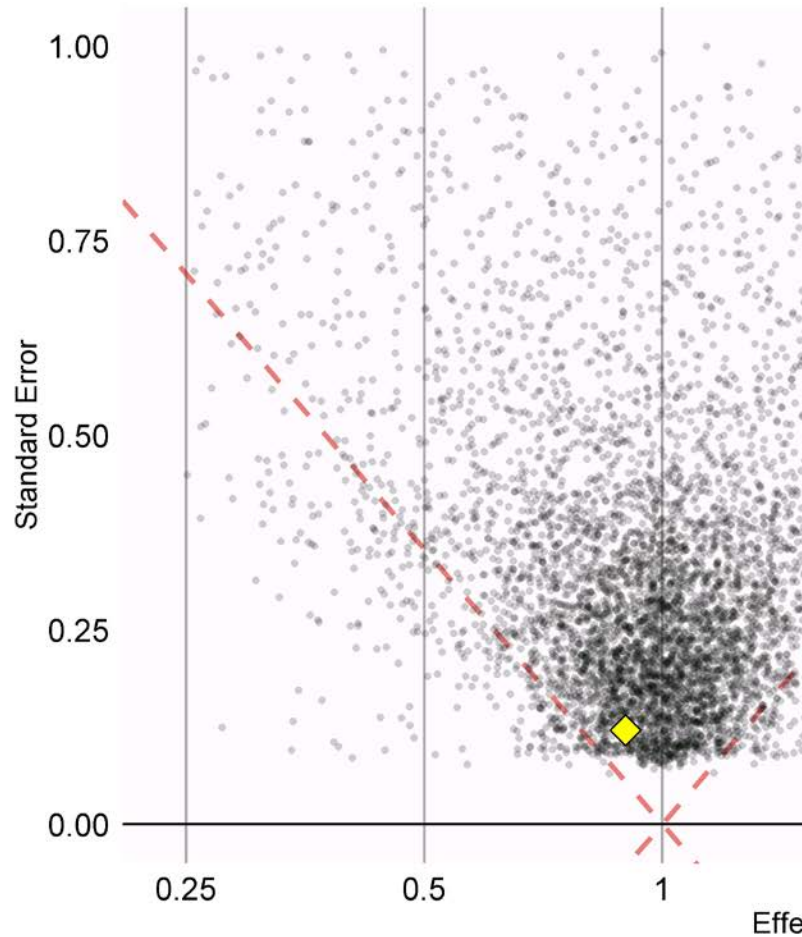
Calibrated HR = 1.05 (0.51 – 2.51)







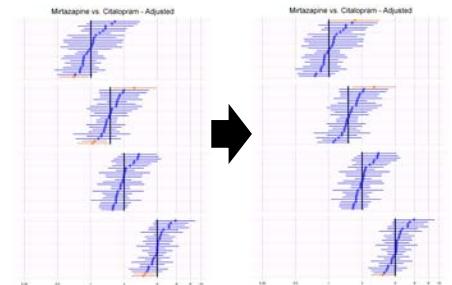
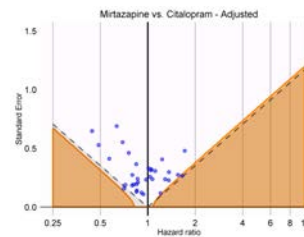
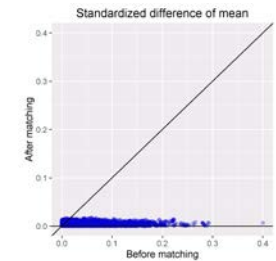
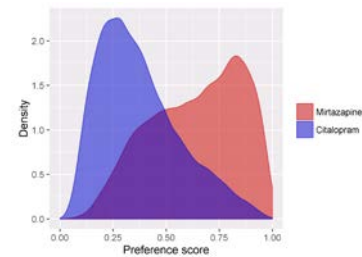
# Example 2



## Mirtazapine vs. Citalopram Constipation

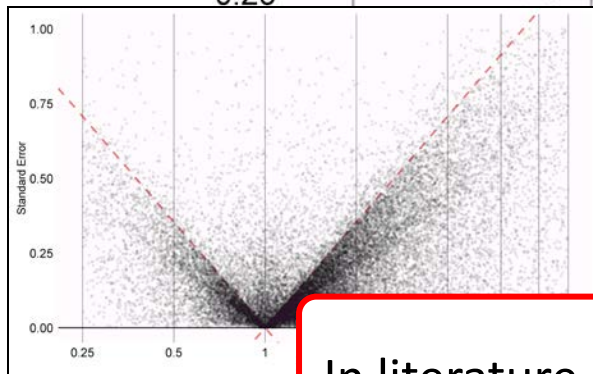
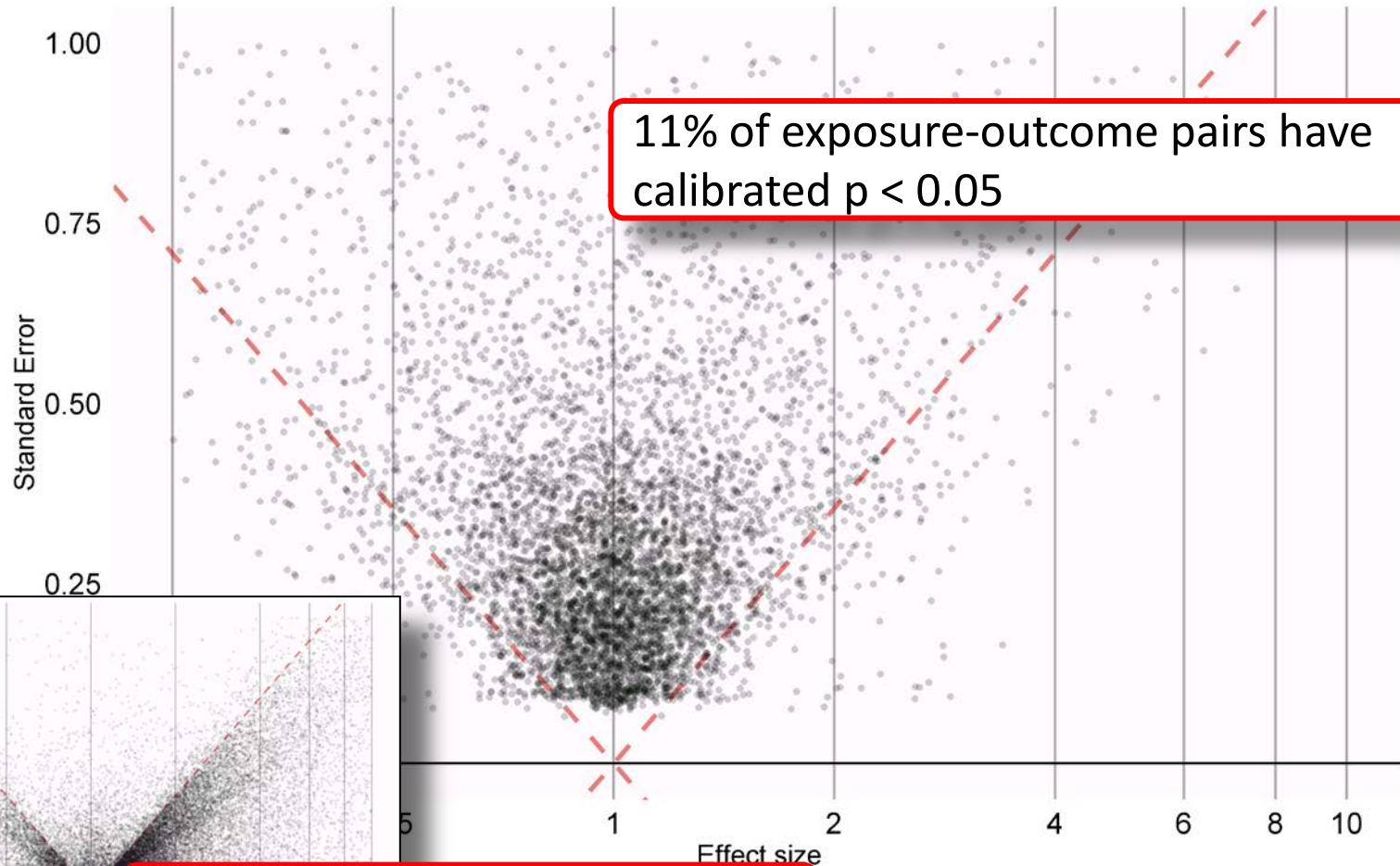
Database: Truven MDCD

Calibrated HR = 0.90 (0.70 – 1.12)





# Estimates are in line with expectations



In literature, 85% have  $p < 0.05$



# Large-scale estimation for depression

- Each estimate produced with same rigor, and could be published as a paper
  - Propensity score adjustment
  - Cox regression
  - Calibrated using negative and positive controls
  - ...
- Not “data-mining”!
  - Results should be interpreted considering multiple testing
  - This can’t be done for literature



# OHDSI recommendations for evidence dissemination

- ✓ Address observation study bias

Addressed by adjusting for confounding, and **verifying** bias was addressed. Disseminate your diagnostics and evaluations.

- ✓ Address publication bias

Avoided by showing all tests that were performed, not just those with  $p < 0.05$

- ✓ Address p-hacking

Very hard to fine-tune analysis to one specific result



# Population-level effect estimation

## Evidence Generation

- Write and share protocol
- Open source study code
- Use validated software
- Replicate across databases

## Evidence Evaluation

- Produce standard diagnostics
- Include negative controls
- Create positive controls
- Calibrate confidence interval and p-value

## Evidence Dissemination

- Don't provide only the effect estimate
- Also share protocol, study code, diagnostics and evaluation
- Produce evidence at scale



# Building the LHC of observational research?



A photograph of several sailboats with white sails on a blue sea under a blue sky with light clouds. The text 'Join the journey' is overlaid in the center.

Join the journey