



Journal of Clinical Epidemiology

Journal of Clinical Epidemiology 187 (2025) 111943

ORIGINAL RESEARCH

Large responses to antidepressants or methodological artifacts? A secondary analysis of STAR*D, a single-arm, open-label, nonindustry antidepressant trial

Colin Xu^a, Florian Naudet^{b,c}, Thomas T. Kim^d, Michael P. Hengartner^e, Mark A. Horowitz^f, Irving Kirsch^g, Joanna Moncrieff^{f,h}, Ed Pigottⁱ, Martin Plöderl^{j,*}

^aDepartment of Psychology and Communication, Institute for Modeling Collaboration and Innovation, University of Idaho, Moscow, ID, USA

^bUniversity of Rennes, CHU Rennes, Inserm, EHESP, Irset (Institut de Recherche en Santé Environnement Et Travail)-UMR_S 1085, CIC 1414 [(Centre d'investigation clinique de Rennes)], Rennes F- 35000, France

^cInstitut Universitaire de France (IUF), Paris, France ^dDepartment of Psychiatry, Weill Cornell Medical College, New York, NY

^eDepartment of Applied Psychology, Kalaidos University of Applied Sciences, Zurich, Switzerland

Research and Development Department, North East London NHS Foundation Trust (NELFT), Rainham, UK

^BHarvard Medical School, Arlington, MA, USA

^BDivision of Psychiatry, University College London, London, UK

Division of Psychiatry, University College London, London, UK

¹Wakefield, RI, USA

^jDepartment of Inpatient Psychotherapy and Crisis Intervention, University Clinic for Psychiatry, Psychotherapy, and Psychosomatics, Christian Doppler Clinic, Paracelsus Medical University, Salzburg, Austria

Accepted 19 August 2025; Published online 25 August 2025

Abstract

Objectives: To replicate Stone et al's (2022) finding that the distribution of response in clinical antidepressant trials is trimodal with large, medium-effect, and small subgroups.

Methods: To apply finite mixture modeling to pre-post Hamilton Depression Rating Scale (HDRS) differences (n = 2184) of STAR*D study's level 1, a single-arm, open-label study. For a successful replication, the best fitting model had to be trimodal, with comparable components as in Stone et al. Secondary/sensitivity analyses repeated the analysis for different baseline levels of depression severity, imputed values, and patient-reported depression symptoms.

Results: The best fitting models were either bimodal or trimodal but the trimodal solution did not meet criteria for replication. The bimodal model had 1 component with HDRS mean change of M = -13.0, SD = 6.7 and included 65.3% of patients, and another component with M = -1.8, SD = 5.1, 34.7%, respectively. For the trimodal model, the component with the largest change (M = -14.3, SD = 6.4) applied to 52% of patients, which differed substantially from the large effect component in Stone et al (M = -18.8, SD = 5.1), which applied to 7.2%. Secondary/sensitivity analyses arrived at similar conclusions, and for patient-reported depression symptoms the best fitting models were unimodal or bimodal.

Conclusion: This analysis failed to identify the trimodal distribution of response reported in Stone et al. In addition to being difficult to operationalize for regulatory purposes, results from mixture modeling are not sufficiently reliable to replace the more robust approach of comparing mean differences in depression rating scale scores between treatment arms. © 2025 The Author(s). Published by Elsevier Inc. This is an open access article under the CC BY license (http://creativecommons.org/licenses/by/4.0/).

Keywords: Antidepressants; Biases; Heterogeneity; Average; Treatment effect; Efficacy

E-mail address: m.ploederl@salk.at (M. Plöderl).

^{*} Corresponding author. Christian Doppler Klinik, Bereich für Stationäre Psychotherapie und Krisenintervention, Ignaz-Harrer-Str. 79, Salzburg 5020, Austria.

What is new?

Key findings

• We could not replicate the trimodal response distribution found in industry trials.

What this adds to what is known?

- Heterogeneity of treatment response can make the average effect unreliable.
- The trimodal response distribution indicates this heterogeneity.
- Modeled response distributions do not seem to be robust and may be prone to biases.

What is the implication and what should change now?

• It is premature to dismiss the average response when evaluating antidepressants.

1. Introduction

The efficacy of antidepressant drugs is usually judged based on statistical significance and the size of the average drug-placebo difference in depression symptom scale scores. Although the average drug-placebo difference is statistically significant, the magnitude is small — about 2 points on the Hamilton Depression Rating Scale (HDRS) [1] which is below nearly all established thresholds of clinically meaningful effects [2,3].

However, efficacy judged by the average drug-placebo difference may be misleading if there is heterogeneity of treatment effects. In 1 study, outcome distributions which were nonnormal and differed between drug-arms and placebo-arms were statistically decomposed into two groups referred to as "nonbenefiters" and "benefiters" and more patients in the drug-arm than in the placeboarm were categorized as "benefiters" [4]. In a much larger recent analysis of individual patient data, the distributions of pre-post changes of depression symptom scores were decomposed into a mixture of different normal subdistributions (modes, components) [1]. It was reported that a trimodal distribution best fit the data overall, which was interpreted to correspond to different subgroups of responders. The three subdistributions corresponded with different levels of improvement, respectively denoted as "large" (with a mean change of M = 16.00 points, standard deviation SD = 4.22), "nonspecific" (M = 8.94, SD = 6.96) and "minimal" (M = 1.68, SD = 2.99) response subgroups (Fig S1). 25% of antidepressanttreated patients were estimated to belong to the

subdistribution of "large" response, compared to 10% of those taking placebo. Stone et al concluded that the small average drug-placebo differences are "best understood as affecting a minority of patients as either an increase in the likelihood of a Large response or a decrease in the likelihood of a Minimal response" (p. 5). While the term "response" is technically appropriate and commonly used, it is problematic as this may suggest that causal processes related to the treatment itself are involved in producing the subdistributions. However, besides the effect of treatment, several other mechanisms are involved in response, such as regression to the mean, natural course or methodological biases. Unfortunately, the findings have been interpreted as showing that there are specific subgroups of participants with distinct "responses", which is misleading since the subdistributions overlap and do not correspond to groups of participants (see examples in the Appendix). In addition, it remains unknown if the trimodal distribution is a robust finding. If such a similar distribution can be found in an open-label, nonindustry study, then this would be compatible with the assumption of a subgroup of patients with a large effect with whatever cause (actual drug effects, natural course, biases, both). If no trimodal distribution is found, especially if there is no subdistribution of a "large effect" then this raises questions about the external validity and interpretation of Stone et al's findings. Therefore, we wanted to explore whether Stone et al's finding of a trimodal distribution could be replicated by applying finite mixture modeling to the STAR*D study, a large single arm, open-label, nonindustry antidepressant trial.

2. Methods

We conducted a secondary analysis of the STAR*D study (level 1). The analysis plan was registered on 2022-11-03 on the Open Science Framework (https://osf.io/rmdu9/) with the protocol uploaded prior to analysis. We used STROBE (strengthening the reporting of observational studies in epidemiology) [5] as a reporting guideline.

2.1. Data

STAR*D is a large publicly funded study [6]. Enrolled patients were 18−75 years of age, seeking care at 18 primary and 23 psychiatric care clinics. Clinical research coordinators screened 4790 patients for major depressive disorder and administered the HDRS, on which 4041 patients scored ≥14, met the other inclusion criteria, and enrolled into the study. In level 1 of the study, all participants were treated open-label with citalopram for up to 14 weeks. HDRS-scores at the end of level 1 were obtained by independent, telephone-based interviewers. The Quick Inventory of Depression Symptomatology (QIDS-SR) self-report rating scale was regularly provided on site. In our main analysis, we selected the 3110 patients who

scored ≥14 on the baseline HDRS because such cutoffs are common in clinical trials, too. In secondary analyses, we included patients independent of their baseline severity. We used a version of the STAR*D data as accessed through the National Institute of Mental Health Data Archive (collection ID #2148) in November 2019 by E.P.

2.2. Analysis

2.2.1. Primary analysis

The primary analysis aimed to replicate the trimodal distribution of pre-post HDRS differences. We used a finite mixture modeling approach where nonnormal distributions are decomposed into a set of different normal distributions, similar to Stone et al [1]. We considered the replication successful if four prespecified criteria were met: a) the best fitting finite mixture model had three components (is trimodal), b) the order of the size of the components is comparable with Stone et al's findings among antidepressant treated patients (ie, the number of "large" responders is smaller than the number of "nonspecific" responders, and the number of "nonspecific" and "minimal" responders differ minimally), c) the proportions of patients in each of the three components are comparable to those found in Stone et al., and d) the mean improvement in the "large" responder group is comparable to that found in Stone et al.

In their publication, Stone et al used a finite mixture model with the data from the drug and placebo groups where means and standard deviations for the drug and placebo groups had to be identical (Fig S1). Because there was no placebo group in the STAR*D trial, we based the comparison on the results provided by Stone for antidepressant-treated patients for a modeling approach where means and standard deviations could vary for the drug and the placebo groups (Fig S1).

We compared the proportions of patients in each of the components found between our study and Stone et al's using χ^2 -tests, based on the 2 x 2 table (study: Stone vs STAR*D, category: large response vs combined unspecific/minimal) and calculated effect sizes. The findings were considered comparable if the upper-bound of the confidence interval of the effect size did not overlap with Cohen's d = 0.3. Similarly, the pre-post differences within the large component were considered comparable if the upper-bound of the effect size did not exceed 0.3.

2.2.2. Secondary analysis

We ran subgroup analyses with all patients independent of their baseline-severity which were categorized into baseline severity HDRS \leq 18, 18 < HDRS \leq 22, and HDRS \geq 22. We also visually compared the distributions of the pre-and post- HDRS scores to those found in Stone et al.

2.2.3. Sensitivity analyses

We ran two sensitivity analyses to see if the findings were sensitive to handling of missing data and choice of the outcome. First, we repeated the main analysis with imputed values for missing HDRS exit scores. Second, we repeated the main analysis with the patient self-report QIDS-SR as the outcome measure. See changes to the protocol for details below (2.4).

2.3. Statistical analysis

We used R (version 4.3.0) [7] for all analyses and the "mixR" package [8] to apply finite mixture modeling. We evaluated model fit using the Bayesian Information Criterion (BIC) and Akaike's Information Criterion (AIC). Effect sizes were calculated using the "esc" package. Imputation was done using the "areg"-function of the "Hmisc" package [9]. Data were first accessed/analyzed between December 2022 and June 2024 (including a break due to lacking resources). Data were curated/provided by T. K. and C.X., and data analysis was performed by M.P. Changes to the protocol are detailed below (2.4). The R-code is available via the Open Science Framework (https://osf.io/rmdu9/).

2.4. Changes to the initial protocol

For the primary analysis, we had originally planned to use the last observation carried forward (LOCF) approach to impute missing values for patients with missing exit HDRS scores. During data analysis, we realized that there were no intermediate assessments of the HDRS and an LOCF approach would lead to an excess of pre- to post-differences of zero, thus obfuscating the finite mixture modeling approach. Therefore, for the primary analysis, we analyzed only those with available exit HDRS scores.

For the secondary analysis, we planned to impute missing exit HDRS values with multiple imputation using variables with less than 90% missing, but without providing further details. We decided (May 2024) to use multiple imputation where, for each imputed sample, "a flexible additive model is fitted on a sample with replacement from the original data and this model is used to predict all of the original missing and nonmissing values for the target variable" [9]. We generated 30 imputed samples, following the recommendation of Harrell (2015) [10]. For imputation we used the baseline demographic and clinical variables (Table 1) and the last available QIDS-SR values because these correlated highly with the HDRS (r = 0.81).

In the protocol we did not prespecify how to calculate Cohen's d for the $\chi 2$ -tests. Because of the theoretical priority to compare the components with the largest pre-post change in HDRS values, we decided (March 2024) to base it on the 2 \times 2 table Study (Stone vs STAR*D) \times category (large response vs combined unspecific/minimal).

Table 1. Sociodemographic and clinical characteristics (primary analysis, N = 3110)

| | M or N | SD or % | % Missing |
|------------------------------------|---------|---------|-----------|
| Demographic features | | | |
| Age | 41.00 | 13.03 | 0 |
| Female | 1996 | 64 | 0 |
| Race | | | |
| Black | 327 | 11 | 0 |
| Hispanic/Latino | 402 | 13 | 0 |
| White | 1977 | 64 | 0 |
| Other | 404 | 13 | 0 |
| Education (y) | 13.60 | 3.23 | 47 |
| Monthly household income | 2321.86 | 2978.19 | 7 |
| Employed | 1711 | 55 | 0 |
| Employed Unemployed | 1711 | 39 | 0 |
| Retired | 170 | 5 | 0 |
| Insurance | 170 | 5 | U |
| Private | 1517 | 4917 | 0 |
| Public | 580 | 19 | 0 |
| None | 1047 | 34 | 0 |
| Marital status | 1047 | 34 | Ŭ |
| Single | 905 | 29 | 0 |
| Married/cohabiting | 1287 | 41 | 0 |
| Separated/divorced | 823 | 26 | 0 |
| Widowed | 92 | 3 | 0 |
| Clinical features | 32 | 3 | O . |
| First episode age < 18 | 1200 | 39 | 1 |
| | 1940 | 67 | 7 |
| Recurrent depression | | | |
| Family history of depression | 1694 | 55 | 2 |
| Age at first episode | 25.14 | 14.29 | 2 |
| Illness duration (y) | 16.11 | 13.47 | 1 |
| Number of episodes | 5.56 | 9.32 | 15 |
| Duration current episode (mo) | 24.88 | 52.02 | 1 |
| Duration current episode $\geq 2y$ | 787 | 26 | 1 |
| QoL questionnaire | 39.07 | 14.26 | 12 |
| SF-12, mental health | 25.58 | 8.06 | 12 |
| SF-12 physical health | 48.61 | 12.13 | 12 |
| Work and social adjustment scale | 24.98 | 8.67 | 12 |
| HDRS-17 | 21.87 | 5.21 | 0 |
| IDS-C30 | 39.07 | 9.64 | 2 |
| QIDS-IVR | 16.88 | 3.31 | 3 |
| Cumulative illness rating scale | | | |
| Categories endorsed | 2.49 | 1.55 | 0 |
| Total score | 4.74 | 3.88 | 0 |
| | | | 10 |
| Severity score | 1.83 | 0.81 | 10 |
| Psychiatric diagnosis screening | 550 | 10 | • |
| Agoraphobia | 559 | 18 | 1 |
| Alcohol abuse/dependency | 371 | 12 | 1 |
| Bulimia | 607 | 20 | 1 |
| Drug abuse/dependency | 234 | 8 | 1 |
| Generalized anxiety disorder | 736 | 24 | 1 |
| Hypochondriasis | 336 | 11 | 1 |

(Continued)

Table 1. Continued

| | M or N | SD or % | % Missing |
|---|--------|---------|-----------|
| OCD | 723 | 23 | 1 |
| Panic disorder | 422 | 14 | 1 |
| PTSD | 387 | 13 | 1 |
| Social anxiety disorder | 963 | 31 | 1 |
| Somatoform disorder | 284 | 9 | 1 |
| Number of axis I comorbid psychiatric disorders | 0.35 | 0.79 | 1 |

IDS-C30, inventory of depressive symptomatology; OCD, obsessive compulsive disorder; PTSD, posttraumatic stress disorder; SF-12, short-form-12 health-survey; QIDS-IVR, interactive voice response; QoL, quality of life.

Variables with less than 10% missing data were used for imputing the missing HDRS exit scores.

3. Results

3.1. Study populations

A flow chart of the patient selection process is provided in Figure 1. Baseline characteristics of patients are provided in Table 1.

3.2. Primary analysis

In the primary analysis using only complete cases, the best fitting models were bimodal according to the BIC and trimodal according to the AIC (Fig 2). However, the trimodal solution did not meet the other three predefined criteria for replication. First, we found a larger proportion of patients in the large response component than the nonspecific response component (52.0% vs 3.2%), whereas Stone et al found the opposite result (7.2% vs 41.8%). Second, the proportions in the three components were different in our study compared to those in Stone et al's. In our study, the proportion in the large response vs all other components was 52% vs 48%,

compared to 7.2% vs 92.8% in Stone et al., resulting in a large difference, $\chi^2(df=1)=5052.5$, P<.01, d=-1.45 (95% CI -1.40 to -1.50). The proportions in the nonspecific and minimal component were about equal in Stone et al but differed substantially in our study, $\chi^2(df=1)=609.7$, P<.01, d=1.34 (95% CI -1.21 to -1.48). Third, the pre-post improvement in the large-response component was -14.3 (95% CI -13.93 to -14.67) in our study and -18.8 (95% CI -18.63 to -18.97) in Stone et al., d=-0.83 (95% CI -0.76 to -0.90). The large response component in our study was more similar to the nonspecific component in Stone et al (M=-14.3 vs M=-14.8).

3.3. Secondary analysis

3.3.1. Results by baseline severity

For the subgroup of patients with a baseline severity HDRS score \leq 18, the best fitting model had two and three components according to the BIC and AIC, respectively (Fig S2). However, the components in the trimodal solution

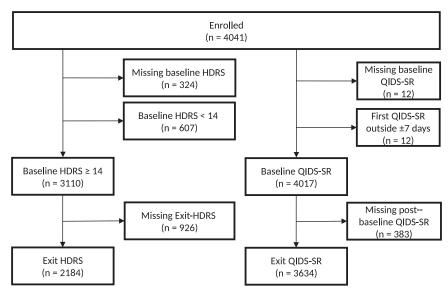
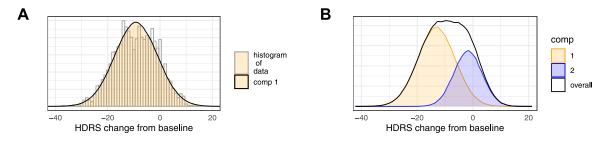
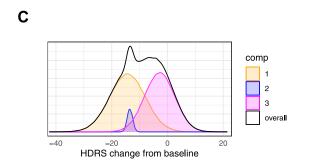


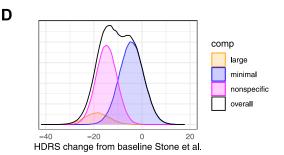
Figure 1. Patient flowchart.



| Model | Component | M (SD) | % | AIC | BIC |
|-------|-----------|--------------|--------|----------|----------|
| 1 | 1 | -9.13 (8.18) | 100.00 | 15384.21 | 15395.59 |

| Model | Component | M (SD) | % | AIC | BIC |
|-------|-----------|---------------|-------|----------|----------|
| 2 | 1 | -13.04 (6.70) | 65.34 | 15351.08 | 15379.53 |
| | 2 | -1.77 (5.07) | 34.66 | | |





| Model | Component | M (SD) | % | AIC | BIC |
|-------|-----------|---------------|-------|----------|----------|
| 3 | 1 | -14.32 (6.41) | 51.99 | 15348.56 | 15394.07 |
| | 2 | -13.47 (1.00) | 3.20 | | |
| | 3 | -2.80 (5.39) | 44.81 | | |

| Model | Component M (SD) | | % | |
|-------|------------------|-------------|------|--|
| 3 | 1 | -18.8 (5.1) | 7.2 | |
| | 2 | -14.8 (4.3) | 41.8 | |
| | 3 | -4.4 (5.1) | 51.0 | |

Figure 2. Results of Finite Mixture Modeling with 1 to 3 components (panels A—C) and, for comparison, the results from the drug-arm in Stone et al (2022) for a model where the means and standard deviation could vary for both arms (panel D). The components are plotted in different colors. The distribution of the original values is plotted as a histogram in the background of panel A. The densities of the mixture models are plotted as thick black lines ("overall"). (For interpretation of the references to color in this figure legend, the reader is referred to the Web version of this article.)

differed from those in Stone et al because the proportions in the small and nonspecific response components differed substantially: 27.3% vs 53.1% in the STAR*D study compared with 51.0% vs 41.8% in Stone et al.

For the subgroup of patients with a baseline severity $18 < HDRS \le 22$, the best fitting model had two components (Fig S3).

For the subgroup of patients with a baseline severity HDRS score >23, the best fitting models had two and 1 component(s) according to the BIC and AIC, respectively (Fig S4). No trimodal solution could be found because the model did not converge.

3.3.2. Visual comparison of distributions of pre- and post-HDRS scores

The baseline HDRS scores of the dataset in Stone et al (their eFigure 2) seemed to be nonnormally distributed with a tighter distribution centered at an HDRS score of 23. In

contrast, the baseline HDRS scores in the STAR*D study seemed to be more normally distributed (Fig S8). We could not compare the posttreatment HDRS scores in the sample as a whole since this information was not available in Stone et al.

3.4. Sensitivity analyses

In the sensitivity analyses where missing exit HDRS scores were imputed, the models converged in 28 of the 30 samples. The best fitting models according to the BIC were unimodal in 3 of the 28 imputed samples, bimodal in 24 samples, and trimodal in 1 sample (Fig S5). According to the AIC, the best fitting model was bimodal in 9 of the 28 imputed samples and trimodal 19 times. The trimodal solutions varied substantially in their nature, meaning that very different solutions fitted the data equally well (online supplement).

For the QIDS-SR as outcome, the best fitting model had 1 and 2 components according to the BIC and AIC, respectively (Fig S7).

4. Discussion

This study analyzed the HDRS pre-post differences in level 1 of the STAR*D study, where all patients were treated with citalogram, to see if the distribution is nonnormal and better explained by subdistributions similar to those in industry-sponsored clinical trials [1]. Using finite mixture modeling, Stone et al [1] found that the nonnormal distribution in their dataset was better explained by a trimodal distribution than a unimodal 1, including a small component with a large mean change from baseline. In contrast, the best fitting model in the STAR*D data was either bimodal or trimodal, but the trimodal model differed substantially from Stone et al's so that none of the prespecified criteria for replication was met. In particular, we did not identify a response component which was comparable with the large response component in Stone et al. We observed similar discrepant findings in secondary analyses for different baseline-levels of depression or for imputed values for the HDRS. In the sensitivity analysis with the self-report QIDS-SR measure, the best fitting model had only one or two components. Thus, none of the results from our analysis of the STAR*D data were in line with the findings by Stone et al.

Our results are relevant for the interpretation of findings from clinical trials, where the average efficacy of antidepressants is small and likely not clinically significant [2,3]. However, the average drug-placebo difference might be misleading if there is heterogeneity in treatment effects. This was suggested by the results of Stone et al's modeling analysis that the outcome distribution is nonnormal and better explained by three response components. Patients classified into the "large" response component were suggested to be "(endo)phenotypes that are specifically responsive to antidepressant drugs (p. 5)". In our study, we could not replicate these findings, that is, finite mixture modeling results did not show that the trimodal model was consistently the best fitting model and there was no comparable distribution of "large" response.

If some patients respond especially well, as suggested by the trimodal distribution in Stone et al., then comparable distributions should also be seen in real-world trials such as STAR*D and not only in clinician rating scales but also in self-report scales. The failure to replicate the trimodal distribution in the STAR*D study and the different findings for clinician vs self-reports raises doubts about the generalizability of the trimodal findings in randomized controlled trials and the finding that there is a subgroup of patients who respond especially well.

How could the discrepant findings be explained? For example, unblinding is present in most trials in which

blinding is tested and this was associated with increased efficacy in some studies [11,12] but not in others [13]. Unblinding may lead to biased symptom ratings where improvement in the drug-arm is overestimated and improvement in the placebo arm is underestimated, leading to a shift of distributions. A tendency of clinicians to harmonize symptom ratings may explain bimodal response distributions, leading to overall scores clustering at either end of the distribution. A trimodal distribution could be explained if the rating bias interacts with degree of symptomreduction, that is, if overestimation of improvement might be stronger for larger symptom-reductions, and underestimation of improvement might be stronger for smaller symptom-reductions. There is some evidence that improvement of symptoms on the HDRS toward symptom remission (eg, from 1 to 0 on an HDRS item) is judged as being more important than other changes (eg, from 3 to 2) [14,15]. If there are several HDRS items rated as zero, then this may bias ratings of other items more strongly toward improvement, compared to scenarios with no or few zero HDRS-item ratings. The effect of these putative biases on outcome distributions could be tested with simulations.

Rating biases may be less pronounced in nonindustry trials such as the STAR*D trial, perhaps explaining why we could not replicate the trimodal distribution in Stone et al. Another finding raising doubt on the robustness of the trimodal solution is that it could not be found for self-report symptom ratings. Here, finite mixture models supported unimodal or bimodal response distributions. It would be interesting to repeat the analysis using selfreport outcomes in industry trials. Furthermore, the distributions of drug and placebo arms should be more similar with more successful blinding. Other explanations for the differences between our results and those by Stone et al are different recruitment and inclusion/exclusion criteria in industry trials and the STAR*D study, the variety of different antidepressants in Stone et al's study, or the use of different strategies to handle missing data. Comparison is also limited because there was no placebo control group in the STAR*D study.

Our results suggest that further research is needed on the distributions of symptom measures in antidepressants trials and consideration of its implications. What constitutes a major deviation from the normal distribution and their difference between drug and placebo has not been discussed adequately yet, to our knowledge. The finite mixture modeling approach may deflect from the small average drug-placebo difference in antidepressant trials. If taken to its logical conclusion then even treatments with zero (or even negative) mean drug-placebo differences cannot be dismissed until subgroups with large improvements can be ruled out, invalidating wellestablished testing paradigms for treatments. Furthermore, there is good reason to remain skeptical about the outlook to identify patients who benefit especially well from treatment [16] and until subgroups of patients who respond more or less to treatment cannot be predicted, results from statistical models to decompose the outcome are just descriptions of data. Unfortunately, the subgroups identified via finite mixture models are easily misinterpreted in at least two ways (examples in the Appendix). First, in the interpretation of Stone et al's study, the overlap between the subdistribution has been ignored and thus the drug-placebo differences in "response" have been overestimated. Second, the subgroups have been interpreted as distinct groups of patients caused by different effects of the treatment but this cannot be inferred from results of statistical models. Finally, reliably identifying subgroups of responders via finite mixture modeling requires large samples and should be seen in patient-reported outcomes, too. For smaller samples, qualitatively different solutions may have comparable fit and the results may be susceptible to imputation methods.

5. Conclusion

In conclusion, the trimodal antidepressant response distribution as reported in Stone et al could not be replicated using data from the STAR*D trial, an open-label, nonindustry sponsored real-world antidepressant study. Therefore, our results do not support the notion that a subgroup of patients with a large response exists. Instead, these findings support the assumption that the putative subgroups from industry randomized controlled trials may be artifacts caused by methodological biases.

CRediT authorship contribution statement

Colin Xu: Writing — review & editing, Writing — original draft, Data curation, Conceptualization. Florian Naudet: Writing - review & editing, Writing - original draft, Conceptualization. Thomas T. Kim: Writing - review & editing, Writing - original draft, Data curation, Conceptualization. Michael P. Hengartner: Writing — review & editing, Writing - original draft, Conceptualization. Mark A. Horowitz: Writing - review & editing, Writing - original draft, Conceptualization. Irving **Kirsch:** Writing – review & editing, Writing – original draft, Conceptualization. Joanna Moncrieff: Writing - review & editing, Writing - original draft, Conceptualization. Ed Pigott: Writing - review & editing, Writing original draft, Conceptualization. Martin Plöderl: Writing - review & editing, Writing - original draft, Supervision, Methodology, Formal analysis, Conceptualization.

Declaration of competing interest

J.M. receives royalties for three books about psychiatric drugs, and is a coapplicant on the REDUCE trial, funded by

the National Institute of Health Research, evaluating digital support for patients stopping long-term antidepressant treatment. M.A.H. and J.M. are both coapplicants on the RELEASE and RELEASE + trials in Australia funded by the National Health and Medical Research Council (NHMRC) and Medical Research Future Fund (MRFF) evaluating hyperbolic tapering of antidepressants. M.A.H. reports being a cofounder of Outro Health which aims to provide digital support for patients in the United States to help stop no longer needed antidepressant treatment using gradual, hyperbolic tapering. There are no competing interests for any other author. F.N. received funding from the French National Research Agency (ANR-17-CE36-0010), the French ministry of health and the French ministry of research. He is a workpackage leader in the OSIRIS project (Open Science to Increase Reproducibility in Science). The OSIRIS project has received funding from the European Union's Horizon Europe research and innovation programme under grant agreement No. 101094725. He is a work-package leader for the doctoral network MSCA-DN SHARE-CTD (HORIZON-MSCA-2022-DN-01 101,120,360), funded by the EU. He is part of the Lown List of Industry-Independent Health Experts.

Acknowledgments

We thank Marc Stone for providing additional analysis and helpful discussions.

Supplementary data

Supplementary data related to this article can be found at https://doi.org/10.1016/j.jclinepi.2025.111943.

Data availability

The authors do not have permission to share data.

References

- [1] Stone MB, Yaseen ZS, Miller BJ, Richardville K, Kalaria SN, Kirsch I. Response to acute monotherapy for major depressive disorder in randomized, placebo controlled trials submitted to the US Food and Drug Administration: individual participant data analysis. BMJ 2022;378:e067606. https://doi.org/10.1136/bmj-2021-067606.
- [2] Hengartner MP, Plöderl M. Estimates of the minimal important difference to evaluate the clinical significance of antidepressants in the acute treatment of moderate-to-severe depression. BMJ Evid Based Med 2021;27:69—73. https://doi.org/10.1136/bmjebm-2020-111600.
- [3] Moncrieff J, Kirsch I. Efficacy of antidepressants in adults. BMJ 2005;331:155-7. https://doi.org/10.1136/bmj.331.7509.155.
- [4] Thase ME, Larsen KG, Kennedy SH. Assessing the "true" effect of active antidepressant therapy ν. placebo in major depressive disorder: use of a mixture model. Br J Psychiatry 2011;199:501–7. https://doi. org/10.1192/bjp.bp.111.093336.
- [5] Von Elm E, Altman DG, Egger M, Pocock SJ, Gøtzsche PC, Vandenbroucke JP. The Strengthening the Reporting of Observational

- Studies in Epidemiology (STROBE) statement: guidelines for reporting observational studies. Int J Surg 2014;12:1495–9. https://doi.org/10.1016/j.ijsu.2014.07.013.
- [6] National Institute of Mental Health. NIMH » sequenced treatment alternatives to relieve depression (STAR*D) Study. NIMH » sequenced treatment alternatives to relieve depression (STAR*D) Study. Available at: https://www.nimh.nih.gov/funding/clinical-research/practical/stard. Accessed September 25, 2022.
- [7] R Core Team. R: A language and environment for statistical computing. Vienna, Austria: R Foundation for Statistical Computing; 2023.
- [8] Yu Y. mixR: Finite mixture modeling for raw and binned data. The Comprehensive R Archive Network; 2021. Available at: https:// cran.r-project.org/web/packages/mixR. Accessed September 10, 2025.
- [9] Harrell F. Hmisc: Harrell miscellaneous. R Package Version 5.1-0. 2023. The Comprehensive R Archive Network. Available at: https://cran.r-project.org/web/packages/Hmisc. Accessed September 10, 2025.
- [10] Harrell FE. Regression modeling strategies: with applications to linear models, logistic and ordinal regression, and survival analysis. Cham: Springer International Publishing; 2015. https://doi. org/10.1007/978-3-319-19425-7.

- [11] Jureidini J, Moncrieff J, Klau J, Aboustate N, Raven M. Treatment guesses in the treatment for adolescents with depression study: Accuracy, unblinding and influence on outcomes. Aust N Z J Psychiatry 2024;58:355–64. https://doi.org/10.1177/00048674231218623.
- [12] Scott AJ, Sharpe L, Colagiuri B. A systematic review and metaanalysis of the success of blinding in antidepressant RCTs. Psychiatry Res 2022;307:114297. https://doi.org/10.1016/j.psychres.2021.114297.
- [13] Lin Y-H, Sahker E, Shinohara K, Horinouchi N, Ito M, Lelliott M, et al. Assessment of blinding in randomized controlled trials of anti-depressants for depressive disorders 2000—2020: a systematic review and meta-analysis. EClinicalMedicine 2022;50:101505. https://doi.org/10.1016/j.eclinm.2022.101505.
- [14] Kim TT, Xu C, DeRubeis RJ. Mapping female patients' judgments of satisfaction to hypothetical changes in depression symptom severity. Behav Ther 2022;53:392–9. https://doi.org/10.1016/j.beth.2021.10.003.
- [15] Kim TT, Xu C, Derubeis RJ. Patients' judgments of the importance of treatment-induced reductions in symptoms of depression: the role of specific symptoms, magnitudes of change, and post-treatment levels. Psychother Res 2022;32:404—13. https://doi.org/10.1080/10503307.2021.1938731.
- [16] Harrell F. Statistical thinking the burden of demonstrating HTE 2019. Available at: https://www.fharrell.com/post/demonte/. Accessed August 29, 2024.