

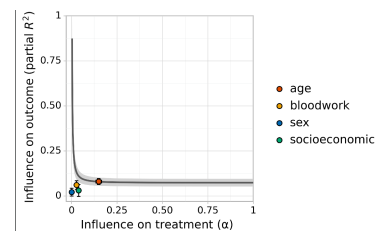
1 We thank the reviewers for the insightful comments and questions, and for their support. Using the camera ready 9th
2 page, we have added a discussion section, and expanded the related work.

3 In the new discussion, we emphasize that Austen plots are most useful in situations where the conclusion from the plot
4 would be ‘obvious’ to a domain expert. For instance, the LaLonde RCT plot shows that a confounder would have to
5 be much stronger than the observed covariates to induce substantial bias. Similarly, the LaLonde observational plot
6 shows that confounders similar to the observed covariates could induce substantial bias. Such conclusions would not be
7 affected by mild perturbations of the dots or the line. By contrast, Austen plots should not be used to draw conclusions
8 such as, “I think a latent confounder could only be 90% as strong as ‘age’, so there is a small non-zero effect”. Such
9 nuanced conclusions might depend on the particular sensitivity model we use, or statistical misestimation or incautious
10 interpretation of the calibration dots—the latter two concerns raised by the reviewers. Drawing precise quantitative
11 (rather than qualitative) conclusions about induced bias from Austen plots would require careful consideration of these
12 issues, and expert statistical guidance. Hence, we recommend that the plots should be used mainly with domain experts
13 to guide qualitative conclusions (“this job program likely works”, “this study doesn’t establish drug efficacy”).

14 **R1 Cinelli and Hazlet (CH)** warn against ‘informal’ benchmarking procedures. Our bias calculation is based on a
15 formal bound, so this critique doesn’t directly apply. To illustrate the issue, CH provide an example where the change
16 in ATE induced by leaving out an observed covariate X is smaller than the bias induced by omitted variable U , even
17 though X and U come into the model in an identical manner. In contrast, the bias estimate in the Austen plot model is
18 typically higher than the change in ATE induced by leaving out a variable—this is shown in our model conservatism
19 experiments. In particular, Austen plots (correctly) anticipate that the bias from omitting U can be higher than the ATE
20 difference induced by computing without X , even when U and X have identical confounding strength.

21 The CH example draws attention to a point that requires some care. The crux of it is 1. strength of influence of U is
22 computed *conditional on the observed covariates*, and 2. the reduction in uncertainty in going from only X to both
23 X, U may be greater than the reduction in uncertainty in going from nothing to X . Thus, it’s possible in principle
24 for the influence of U to be larger than the estimated (observed) influence of X , even if X and U are similar causal
25 variables. In other words, the omission of U could cause us to underestimate the influence of X from the observed
26 data. Accordingly, when the domain expert compares the strength of the unobserved confounder to a reference dot
27 for X on the plot, they must also ask if knowing U could have substantially increased the predictive power of X . In
28 cases where this seems plausible—e.g., the domain expert expects an important interaction between U and X —then
29 naively eyeballing the dot vs line position may be unreliable and further careful thought is required. However, we note
30 that examples of this kind are somewhat contrived. Indeed, we usually expect the opposite effect. If U and X are
31 dependent, then some of the information in X will be redundant, and the measured R^2 and α will *overestimate* the true
32 influence. That is, the CH effect tends to make our sensitivity analysis conservative. This is reflected by the fact that
33 grouping similar covariates (to mitigate redundancy) led to higher computed influence in every example considered
34 in the paper. Although the CH effect is an important conceptual point and the domain expert should consider it as
35 part of due diligence, it doesn’t seem to have much impact on the practical use of Austen plots. We also note that
36 this conceptual subtlety is a generic feature of calibrating sensitivity analyses, not particular to our method. This is
37 reflected, e.g., in Franks et al §5.2.1, which describes their calibration procedure based on looking at variance explained
38 by observed Z conditioned on $X \setminus Z$ —a procedure similar to ours, carrying the same nuance. We have clarified this
39 point in the newly added discussion section. We thank the reviewer for bringing this to our attention.

40 **R2 Inference** We agree that inference is a key issue. We have added some addi-
41 tional detail estimation via the plug-in estimators. For handling uncertainty, we sug-
42 gest bootstrapping several plot versions and trusting only conclusions supported by all plots.
43 We tried this on the examples in the paper, and found no change in conclusions. We have
44 added an appendix describing how to visually summarize the uncertainty; see fig (cf. Fig 1 in
45 paper). Note: α uncertainty is plotted, but too small to show up clearly. Formal uncertainty
46 quantification is an interesting direction for future work; particularly important because
47 bootstrapping requires model refitting, which can be computationally intensive for the ML
48 models that motivate the paper.



49 We emphasize that the unobserved confounder bias we address exists even in the (very)
50 large data regime where plug-in estimators work very well. Indeed, this is the setting where
51 misidentification bias matters most relative to statistical error. Accordingly, we believe the paper, using the plug-in
52 estimators, is a significant contribution even in the absence of efficiency guarantees. In the discussion, we caution
53 that the plots may be misleading if there’s large uncertainty in the estimator of Q and g . As part of the discussion on
54 estimation issues, we note efficient estimation as a good direction for future work.

55 **Related work** Thank you for the pointers to the additional related literature. As you suggest, we have used the extra
56 page to substantially expand the related work section.