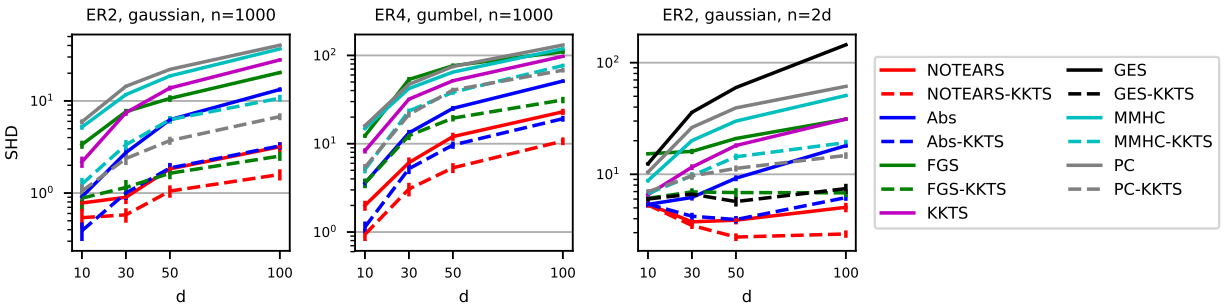


1 We thank the reviewers for their efforts. We are glad that R1, R2, R3 appreciated the improved understanding of  
 2 NOTEARS (“useful negative results”). Thanks to R2 for describing our theoretical insights in general as “a potentially  
 3 useful side effect for others” and “surprisingly readable for a paper with lots of theoretical results.” R1, R2, R4 recognize  
 4 our empirical results to be “successful” and “a significant improvement.” Below we respond to reviewer comments.

5 **R3: Questionable that NOTEARS, FGS outperform earlier methods, [2, Table 1] shows MMHC, PC perform**  
 6 **much better than GES (here FGS) for moderate  $d$ .** Thank you for pointing this out. We agree that the statement  
 7 on line 276 is too broad and will remove it. To address R3’s concern, we first compared with MMHC and PC in the  
 8 experimental setting of Sec. 5. The significance level  $\alpha$  was chosen from the range considered in [2] to minimize SHD.  
 9 The two left panels below show that while MMHC and PC do not perform better than FGS, they are also significantly  
 10 improved by KKTS (we will report full results in the paper). We then performed a second experiment with  $n = 2d$   
 11 to be closer to the setting of [2, Table 1], also adding the GES implementation used in [2]. The right panel shows that FGS  
 12 is actually an improvement over GES, remaining better than MMHC and PC (except  $d = 10$ ) while GES is worse.



13 **R3: “Paper is fairly incremental, developing a single heuristic local search method (namely NOTEARS that**  
 14 **enjoys no non-trivial performance guarantees).”** We naturally disagree about incrementality and are not sure what is  
 15 meant by non-trivial. Prop. 3 provides a negative guarantee for NOTEARS (which is *not* our method), whereas Thms 9  
 16 and 7 provide positive guarantees for our KKTS method to yield KKT points and local minima. We would thus not call  
 17 KKTS heuristic or lacking guarantees. To get from Prop. 3 to KKTS requires several more contributions: reformulating  
 18 the problem and proving that KKT conditions are necessary (Thm 6), and relating the KKT conditions to edge absence  
 19 constraints (Lem. 6, Thm 8). Sec. 2 makes additional contributions in generalizing acyclicity constraints from [32,30].

20 **R4: Apparently only applicable to continuous case, no mention of categorical. More needed on limitations and**  
 21 **position in literature.** We thank R4 for prompting us to elaborate upon the problem setting and what remains for future  
 22 work, which we will do in the paper. The theory and methods apply straightforwardly to binary variables (although  
 23 we have not experimented with them) but not non-binary categorical variables. As Sec. 3, para. 1 states, the key  
 24 assumption is that each edge is associated with a single parameter  $W_{ij}$ . In a (generalized) linear structural equation, a  
 25 single parameter can account for the effect of a binary or continuous input, while a binary output can be handled by a  
 26 suitable loss function (e.g. logistic). However, a single parameter is likely insufficient for a non-binary categorical input  
 27 (typically encoded into multiple binary variables) or output (e.g. [14] proposes multi-logit regression with parameters  
 28 for each level). Therefore the extension to multiple parameters per edge (Sec. 6) is desirable to address categorical  
 29 variables as well as nonlinear models. **Abstract:** We will add a sentence on the one-parameter-per-edge assumption.  
 30 **Title:** We find it difficult to capture this assumption in a few readily understood words, but perhaps R4 has a suggestion.

31 **R1: “What leads to better or worse SHD...F (always squared error...danger of overfitting?), thresholding,**  
 32 **acyclicity constraint and even centering.”** Thanks for the thoughtful questions. Below we summarize what we  
 33 know/have reported. We think a proper exploration would best be left to a journal extension of this paper. **Score**  
 34 **function  $F$ :** By keeping this as least squares, we have somewhat avoided overfitting to the noise type (Gaussian,  
 35 Gumbel, etc., usually unknown) as opposed to using the log-likelihood. [32, Sec. 5.3] shows that NOTEARS can  
 36 achieve scores close to those of the exact optimizer GOBNILP (especially before thresholding), i.e. it is fitting well  
 37 but not over-fitting, but we have not done a similar comparison for NOTEARS-KKTS. **Acyclicity constraint:** The  
 38 NOTEARS vs. Abs comparison shows that setting  $A = W \circ W$  is empirically superior to  $A = |W|$ . As for the function  
 39  $h$ , Appendix C.1 states that the polynomial  $h(A)$  from [30] performs slightly better and is slightly faster than the  
 40 exponential  $h(A)$  from [32] (the authors of [32] seem to agree as their code now uses polynomial  $h$ ). Section 6 mentions  
 41 that other  $h$  in the class of eq. (1) could be explored in future work. **Thresholding:** We followed [32] and fixed the  
 42 threshold  $\omega = 0.3$  for NOTEARS and for our NOTEARS-inspired algorithms (Abs, KKTS), in part to demonstrate  
 43 success without too much parameter tuning. But we agree that the role of thresholding could be further explored.  
 44 **Centering:** We were also surprised by the effect this had, as reported in Appendix C.3.

45 **R1: “Fully continuous” is overstatement.** We will remove “fully”. Note that NOTEARS [32] also uses thresholding.