
Beware of the Simulated DAG!

Causal Discovery Benchmarks May Be Easy To Game

Alexander G. Reisach^{1,2}

Christof Seiler^{2,3}

Sebastian Weichwald¹

¹Department of Mathematical Sciences, University of Copenhagen, Denmark

²Department of Data Science and Knowledge Engineering, Maastricht University, The Netherlands

³Mathematics Centre Maastricht, Maastricht University, The Netherlands

Abstract

Simulated DAG models may exhibit properties that, perhaps inadvertently, render their structure identifiable and unexpectedly affect structure learning algorithms. Here, we show that marginal variance tends to increase along the causal order for generically sampled additive noise models. We introduce *varsortability* as a measure of the agreement between the order of increasing marginal variance and the causal order. For commonly sampled graphs and model parameters, we show that the remarkable performance of some continuous structure learning algorithms can be explained by high *varsortability* and matched by a simple baseline method. Yet, this performance may not transfer to real-world data where *varsortability* may be moderate or dependent on the choice of measurement scales. On standardized data, the same algorithms fail to identify the ground-truth DAG or its Markov equivalence class. While standardization removes the pattern in marginal variance, we show that data generating processes that incur high *varsortability* also leave a distinct covariance pattern that may be exploited even after standardization. Our findings challenge the significance of generic benchmarks with independently drawn parameters. The code is available at <https://github.com/Scriddie/Varsortability>.

1 Introduction

Causal structure learning aims to infer a causal model from data. Academic disciplines anywhere from biology, medicine, finance, to machine learning are interested in causal models [Rothman et al., 2008, Imbens and Rubin, 2015, Sanford and Moosa, 2012, Schölkopf, 2019]. Causal models not only describe the observational joint distribution of variables but also formalize predictions under interventions and counterfactuals [Spirtes et al., 2000, Pearl, 2009, Peters et al., 2017]. Directed acyclic graphs (DAGs) are common to represent causal structure: nodes represent variables and directed edges point from cause to effect representing the causal relationships. This graphical representation rests on assumptions which have been critically questioned, for example by Dawid [2010]. Inferring causal structure from observational data is difficult: Often we can only identify the DAG up to its Markov equivalence class (MEC) and finding high-scoring DAGs is NP-hard [Chickering, 1996, Chickering et al., 2004]. Here, we focus on learning the DAG of linear additive noise models (ANM).

Data scale and marginal variance may carry information about the data generating process. This information can dominate benchmarking results, such as, for example, the outcome of the NeurIPS Causality 4 Climate competition [Runge et al., 2020]. Here, the magnitude of regression coefficients was informative about the existence of causal links such that ordinary regression-based methods on raw data outperformed causal discovery algorithms [Weichwald et al., 2020]. Multiple prior works

state the importance of data scale for structure learning either implicitly or explicitly. Structure identification by *ICA-LiNGAM* [Shimizu et al., 2006], for example, is susceptible to rescaling of the variables. This motivated the development of *DirectLiNGAM* [Shimizu et al., 2011], a scale-invariant causal discovery algorithm for linear non-Gaussian models. The causal structure of ANMs is proven to be identifiable given the noise scale (cf. Section 2.2). Yet, such identifiability results require knowledge about the ground-truth data scale.

Simulated DAGs may be identifiable from marginal variances under generic parameter distributions. An instructive example is the causal graph $A \rightarrow B$ with structural equations $A = N_A$ and $B = wA + N_B$ with $w \neq 0$ and independent zero-centered noise variables N_A, N_B . The mean squared error (MSE) of a model $X \rightarrow Y$ is given by $\text{MSE}(X \rightarrow Y) = \text{Var}(X) + \text{Var}(Y|X)$. It holds that $\text{MSE}(A \rightarrow B) < \text{MSE}(B \rightarrow A) \iff \text{Var}(A) < \text{Var}(B) \iff (1 - w^2) \text{Var}(N_A) < \text{Var}(N_B)$ (see Appendix A). Deciding the directionality of the edge between A and B based on the MSE amounts to inferring an edge from the lower-variance variable to the higher-variance variable. For error variances $\text{Var}(N_A) \leq \text{Var}(N_B)$ and any non-zero edge weight w , the MSE-based inference is correct. This resembles known scale-based identifiability results based on equal or monotonically increasing error variances [Peters and Bühlmann, 2014, Park, 2020]. However, if the observations of A were multiplied by a sufficiently large constant, the MSE-based inference would wrongfully conclude that $A \leftarrow B$. This is problematic since simply choosing our units of measurement differently may change the scale and variance of A . Arguably, this is often the case for observations from real-world systems: There is no canonical choice as to whether we should pick meters or yards for distances, gram or kilogram for weights, or yuan or dollar as currency. A researcher cannot rely on obtaining the same results for different measurement scales or after re-scaling the data when applying any method that leverages the data scale (examples include Peters and Bühlmann [2014], Park [2020], or Zheng et al. [2018], who employ the least squares loss studied by Loh and Bühlmann [2014]).

Continuous causal structure learning algorithms optimize model fit under a differentiable acyclicity constraint [Zheng et al., 2018]. This allows for the use of continuous optimization and avoids the explicit combinatorial traversal of possible causal structures. This idea has found numerous applications and extensions [Lachapelle et al., 2019, Lee et al., 2019, Ng et al., 2020, Yu et al., 2019, Brouillard et al., 2020, Pamfil et al., 2020, Wei et al., 2020, Zheng et al., 2020, Bhattacharya et al., 2021]; Vowels et al. [2021] provide a review. *NOTEARS* [Zheng et al., 2018] uses the MSE with reference to Loh and Bühlmann [2014], while *GOLEM* [Ng et al., 2020] assesses model fit by the penalized likelihood assuming a jointly Gaussian model. On simulated data and across noise distributions, both methods recover graphs that are remarkably close to the ground-truth causal graph in structural intervention distance (SID) and structural hamming distance (SHD). We agree with the original authors that these empirical findings, especially under model misspecification and given the non-convex loss landscape, may seem surprising at first. Here, we investigate the performance under data standardization and explain how the causal order is (partially) identifiable from the raw data scale alone in common generically simulated benchmarking data.

Contribution. We show that causal structure drives the marginal variances of nodes in an ANM and can lead to (partial) identifiability. The pattern in marginal variances is dominant in ANM benchmark simulations with edge coefficients drawn identically and independently. We introduce varsortability as a measure of the information the data scale carries about the causal structure. We argue that high varsortability affects the optimization procedures of continuous structure learning algorithms. Our experiments demonstrate that varsortability dominates the optimization and helps achieve state-of-the-art performance provided the ground-truth data scale. Data standardization or an unknown data scale remove this information and the same algorithms fail to recover the ground-truth DAG. Even methods using a score-equivalent likelihood criterion (*GOLEM*) recover neither ground-truth DAG nor its MEC on standardized data. To illustrate that recent benchmark results depend heavily on high varsortability, we provide a simple baseline method that exploits increasing marginal variances to achieve state-of-the-art results on these benchmarks. We thereby provide an explanation for the unexpected performance of recent continuous structure learning algorithms in identifying the true DAG. Neither algorithm dominates on raw or standardized observations of the analyzed real-world data. We show how, even if data is standardized and even in non-linear ANMs, a causal discovery benchmark may be gamed due to covariance patterns. Consequently, recent benchmark results may not transfer to (real-world) settings where the correct data scale is unknown or where edge weights are not drawn independent and identically distributed (iid). We conclude that structure learning

benchmarks on ANMs with generically sampled parameters may be distorted due to unexpected and perhaps unintended regularity patterns in the data.

2 Background

2.1 Model Class

We consider acyclic linear additive noise models. Single observations are denoted by $x^{(i)} \in \mathbb{R}^d$ where $x_j^{(i)}$ denotes the j^{th} dimension of the i^{th} iid observation of random vector $X = [X_1, \dots, X_d]^{\top}$. All observations are stacked as $\mathbf{X} = [x^{(1)}, \dots, x^{(n)}]^{\top} \in \mathbb{R}^{n \times d}$ and $x_j \in \mathbb{R}^n$ refers to the j^{th} column of \mathbf{X} . Analogously, $n^{(i)}$ denotes the corresponding i^{th} iid observation of the random noise variable $N = [N_1, \dots, N_d]^{\top}$ with independent zero-centred components. The linear effect of variable X_k on X_j is denoted by $w_{k \rightarrow j} = w_{kj}$. The causal structure corresponding to the adjacency matrix $W = [w_{kj}]_{k,j=1,\dots,d}$ with columns $w_j = [w_{k \rightarrow j}]_{k=1,\dots,d} \in \mathbb{R}^d$ can be represented by a directed acyclic graph $G = (V_G, E_G)$ with vertices $V_G = \{1, \dots, d\}$ and edges $E_G = \{(k, j) : w_{k \rightarrow j} \neq 0\}$. Edges can be represented by an adjacency matrix E such that the $(k, j)^{\text{th}}$ entry of E^l is non-zero if and only if a directed path of length l from k to j exists in G . For a given graph, the parents of j are denoted by $\text{PA}(j)$. The structural causal model is $X = W^{\top} X + N$.

2.2 Identifiability of Additive Noise Models

Identifiability of the causal structure or its MEC requires causal assumptions. Under causal faithfulness and Markov assumptions, the causal graph can be recovered up to its MEC [Chickering, 1995, Spirtes et al., 2000]. Faithfulness, however, is untestable [Zhang and Spirtes, 2008]. Shimizu et al. [2006] show that under the assumptions of no unobserved confounders, faithfulness, linearity, and non-Gaussian additive noise, the causal graph can be recovered from data. Hoyer et al. [2009] show that this holds for any noise distribution under the assumption of strict non-linearity. This finding is generalized to post-nonlinear functions by Zhang and Hyvarinen [2009]. Peters and Bühlmann [2014] prove that the causal structure of a linear causal model with Gaussian noise is identifiable if the error variances are equal or known. Any unknown re-scaling of the data breaks this condition. For the case of linear structural causal models, Loh and Bühlmann [2014] provide a framework for DAG estimation based on a noise variance-weighted least squares score function. For ANMs, they give conditions under which the general Gaussian case can be identified via approximating it by the equal noise-variance case given knowledge of the (approximate) noise scale. Finally, subsuming further prior results on (linear) ANMs [Hoyer et al., 2009, Ghoshal and Honorio, 2017, 2018, Chen et al., 2019], Park [2020] shows that the causal structure is identifiable under regularity conditions on the conditional variances along the causal order. In particular, identifiability holds if the error variances of nodes are weakly monotonically increasing along the causal order.

2.3 Structure Learning Algorithms

Combinatorial structure learning algorithms (such as PC, FGES, DirectLiNGAM) separately solve the combinatorial problem of searching over structures and finding the optimal parameters for each structure. To remain computationally feasible, the search space of potential structures is often restricted or traversed according to a heuristic. One can, for example, carefully choose which conditional independence statements to evaluate in constraint-based algorithms, or employ greedy (equivalence class) search in score-based algorithms. In our experiments, we consider PC [Spirtes and Glymour, 1991], FGES [Meek, 1997, Chickering, 2002b], DirectLiNGAM [Shimizu et al., 2011], and a greedy DAG search (GDS) algorithm MSE-GDS that greedily includes those edges that reduce the MSE the most. For details see Appendix D.

Continuous structure learning algorithms (such as NOTEARS and GOLEM) employ continuous optimization to simultaneously optimize over structures and parameters. As a first step towards expressing causal structure learning as a continuous optimization problem, Aragam and Zhou [2015] propose l^1 -regularization instead of the conventional l^0 -penalty for model selection. Zheng et al. [2018] propose a differentiable acyclicity constraint, allowing for end-to-end optimization of score functions over graph adjacency matrices. We examine and compare the continuous structure learning algorithms NOTEARS [Zheng et al., 2018] and GOLEM [Ng et al., 2020]. For details see Appendix D.

NOTEARS [Zheng et al., 2018] minimizes the MSE between observations and model predictions subject to a hard acyclicity constraint. The MSE with respect to W on observations \mathbf{X} is defined as $\text{MSE}_{\mathbf{X}}(W) = \frac{1}{n} \|\mathbf{X} - \mathbf{X}W\|_2^2$ where $\|\cdot\|_2 = \|\cdot\|_F$ denotes the Frobenius norm.

GOLEM [Ng et al., 2020] performs maximum likelihood estimation (MLE) under the assumption of a Gaussian distribution with equal (EV) or non-equal (NV) noise variances. There are soft acyclicity and sparsity constraints. The unnormalized negative likelihood-parts of the objective function are $\mathcal{L}_{EV}(W, \mathbf{X}) = \log(\text{MSE}_{\mathbf{X}}(W))$ and $\mathcal{L}_{NV}(W, \mathbf{X}) = \sum_{j=1}^d \log(\frac{1}{n} \|x_j - \mathbf{X}w_j\|_2^2)$, respectively, omitting a $-\log(|\det(I - W)|)$ term that vanishes when W represents a DAG [Ng et al., 2020].

To ease notation, we sometimes drop the explicit reference to \mathbf{X} when referring to MSE, \mathcal{L}_{EV} , \mathcal{L}_{NV} .

3 Varsortability

The data generating process may leave information about the causal order in the data scale. We introduce varsortability as a measure of such information. When varsortability is maximal, the causal order is identifiable. Varsortability is high in common simulation schemes used for benchmarking causal structure learning algorithms. We describe how continuous structure learning algorithms are affected by marginal variances and how they may leverage high varsortability. This elucidates the results of continuous methods reported by Zheng et al. [2018], Ng et al. [2020], and others on raw data and predicts impaired performance on standardized data as confirmed in Section 4. We introduce *sortnregress* as simple baseline method that sorts variables by marginal variance followed by parent selection. The performance of *sortnregress* reflects the degree of varsortability in a given setting and establishes a reference baseline to benchmark structure learning algorithms against.

3.1 Definition of Varsortability

We propose varsortability as a measure of agreement between the order of increasing marginal variance and the causal order. For any causal model over variables $\{X_1, \dots, X_d\}$ with (non-degenerate) DAG adjacency matrix E we define varsortability as the fraction of directed paths that start from a node with strictly lower variance than the node they end in, that is,

$$v := \frac{\sum_{k=1}^{d-1} \sum_{i \rightarrow j \in E^k} \text{increasing}(\text{Var}(X_i), \text{Var}(X_j))}{\sum_{k=1}^{d-1} \sum_{i \rightarrow j \in E^k} 1} \in [0, 1] \text{ where } \text{increasing}(a, b) = \begin{cases} 1 & a < b \\ 1/2 & a = b \\ 0 & a > b \end{cases}$$

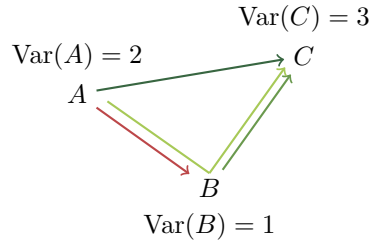
For example, we calculate the varsortability as $v = \frac{1+1+1}{1+1+1+1} = \frac{3}{4}$ given the causal graph below.

Varsortability equals one if the marginal variance of each node is strictly greater than that of its causal ancestors. Varsortability equals zero if the marginal variance of each node is strictly greater than that of its descendants. Varsortability does not depend on choosing one of the possibly multiple causal orders and captures the overall agreement between the partial order induced by the marginal variances and all pathwise descendant relations implied by the causal structure. In the two-node introductory example (cf. Section 1), varsortability $v = 1$ is equivalent to

$$\text{Var}(A) < \text{Var}(B) \iff \text{Var}(N_A) < w^2 \text{Var}(N_A) + \text{Var}(N_B)$$

where A and B are nodes in the causal graph $A \xrightarrow{w} B$ with noise variances $\text{Var}(N_A)$ and $\text{Var}(N_B)$.

We can also understand varsortability as a property of the distribution of graphs and parameters that we sample from for benchmarks on synthetic data. The distribution of weights and noise variances determines whether the causal order of any two connected nodes in the graph agrees with the order of increasing marginal variance and in turn determines the varsortability of the simulated causal models. We observe that even for modest probabilities of any two neighboring nodes being correctly ordered by their marginal variances, the variance of connected nodes tends to increase quickly along the causal order for many ANM instantiations (cf. Appendix G.3). For a heuristic explanation, recall that we obtain the marginal variance of a node by adding the variance contribution of all its ancestors to the



node’s own noise variance; to obtain the variance contribution of an ancestor, we take the product of the edge weights along each directed path from ancestor to node, sum these path coefficient products, square, and multiply with the ancestor’s noise variance. While the sum of path coefficient products may vanish or be small such that the variance contribution of an ancestor cancels out or is damped across the different connecting paths, we find it is unlikely if edge weights are drawn independently (cf. Meek [1995] for why exact cancellation and faithfulness violations are unlikely). Furthermore, the further apart a connected pair of nodes, the more variance may be accumulated in the descendant node along all incoming paths in addition to one ancestor’s (possibly damped) variance contribution further fueling the tendency for descendant nodes to have higher variance than their ancestors. In practice we indeed find that an ordering of nodes by increasing marginal variance closely aligns with the causal order for commonly simulated linear and non-linear ANMs (cf. Appendix G).

3.2 Varsortability and Identifiability

If varsortability $\nu = 1$, the causal structure is identifiable. It can be recovered by ordering the nodes by increasing marginal variance and regressing each node onto its predecessors using conventional sparse regression approaches. The causal structure learning problem is commonly decoupled into causal order estimation and parent selection [Shimizu et al., 2006, Shojaie and Michailidis, 2010, Bühlmann et al., 2014, Chen et al., 2019, Park, 2020]. This decoupling is further warranted, since we only need the causal ordering to consistently estimate interventional distributions [Bühlmann et al., 2014, Section 2.6]. At $\nu = 1$, an ordering by marginal variance is a valid causal ordering. Given a causal ordering, one can construct a fully connected DAG and use parent selection to prune edges and reconstruct the graph in the sample limit under mild assumptions. Bühlmann et al. [2014, Section 2.5] discuss parent selection and Shojaie and Michailidis [2010] establish the consistency of an adaptive lasso approach for edge selection given a valid causal order. The identifiability conditions by Park [2020] are closely related to varsortability, though not equivalent as we prove in Appendix C. Identifiability of the causal order is immediate if varsortability $\nu = 1$, though, this shares severe drawbacks with other identifiability conditions that rely on data scale by Peters and Bühlmann [2014], Loh and Bühlmann [2014], and Park [2020]. First, it is difficult to verify or assess the plausibility of assumptions about the correctness or suitability of the data scale for any given dataset. Second, any unknown rescaling may break previously met identifiability conditions. Third, even if variables are on the same measurement scale the units may not correctly capture the ground-truth causal scale. For example, a dartist’s distance from the dartboard may affect the precision of their throw measured by the distance between hit and target. Here, the effect variable’s marginal variance may be smaller than that of the cause (even) if both distances are measured in centimetres. Nonetheless, it may be possible to exploit varsortability if one can establish that certain assumptions on the data scale be met.

3.3 Varsortability in Benchmarking Scenarios

For real-world data we cannot readily assess nor presume varsortability as we do not know the parameters and data scale of the data generating process. When benchmarking causal structure learning algorithms, however, we can evaluate varsortability for the simulated DAGs and parameter settings. We may acquire an intuition about the probability of varsortable cause-effect pairs in our simulation settings by considering two neighboring nodes $A \xrightarrow{w} B$ in the sampled graph without common ancestors and no other directed path from A to B . Under these assumptions, $\text{Var}(B) = w^2 \text{Var}(A) + \text{Var}(\sum_{C \in \text{PA}(B) \setminus \{A\}} w_{C \rightarrow B}^2 C) + \text{Var}(N_B)$ such that $|w| > 1$ implies that the variable pair is varsortable. To simulate an ANM, we need to sample a DAG, decide on the noise distributions, and then sample edge weights and noise variances. Across simulation instances in a given benchmark set-up, the edge weights $w_{k \rightarrow j}$ and noise variances s_k^2 are iid instantiations of independent random variables W and S^2 . The distributions of W and S^2 induce a distribution of the marginal variance V_Y of node Y in the resulting ANM. The probability for the variable pair $A \rightarrow B$ to be varsortable in a simulated ANM is then bounded from below by $\text{P}[(1 - W_{A \rightarrow B}^2)V_A < S_{N_B}^2]$ (cf. Appendix B). If A is a root node, $V_A = S_{N_A}^2$. In the experimental settings used by, for example, Zheng et al. [2018, 2020], Lachapelle et al. [2019], Ng et al. [2020], edge weights are independently drawn from a uniform distribution and noise standard deviations or variances are either fixed or also drawn independently from a uniform distribution. For our parameters $W \stackrel{\text{iid}}{\sim} \text{Unif}((-2, -0.5) \cup (0.5, 2))$ and $S \stackrel{\text{iid}}{\sim} \text{Unif}((0.5, 2))$, which resemble common choices in the literature, any pair is varsortable with probability at least $2/3$ due to $\text{P}[|W| > 1] = 2/3$, and with probability $p > 0.93$ provided A is a

root node. Empirically, we find that varsortability averages above 0.94 in our simulated graphs and above 0.71 in commonly considered non-linear ANMs (cf. Appendix G). This result indicates that in benchmark simulations the marginal variance of any two nodes in a graph tends to increase along the causal order and that we may game these benchmarks and perform well by exploiting this pattern.

If A and B have a common ancestor or mediator C , the effect of C on B may either compound or partially cancel out the effect of A on B . In practice, the effect commonly increases the variance of the effect node B , which may be attributed to the independent sampling of path coefficients which also renders faithfulness violations improbable [Meek, 1995]. We find varsortability to increase with graph density and the lower bound presented above to be loose. Motivated by the strong impact of different levels of varsortability on some structure learning algorithms as reported in Section 4 and Appendix H.2, we advocate an empirical evaluation and reporting of varsortability (cf. Appendix G.4 for the implementation) when simulating ANMs. We emphasize that even for varsortability < 1 , where the order of increasing variance does not perfectly agree with the causal order, experimental results may still be largely driven by the overall agreement between increasing marginal variance and causal order. The extent to which varsortability may distort experimental comparisons of structure learning algorithms on linear ANMs is demonstrated in Section 4.

3.4 Marginal Variance yields Asymmetric Gradients for Causal and Anti-Causal Edges

We explain how varsortability may dominate the performance of continuous structure learning algorithms. We do not expect combinatorial structure learning algorithms that use a score-equivalent (see e.g. Yang and Chang [2002], Chickering [2002a]) criterion or scale-independent (conditional) independence tests to be dependent on the data scale. This includes PC , as local constraint-based algorithm, $FGES$ as locally greedy score-based search using a score-equivalent criterion, and $DirectLiNGAM$, a procedure minimizing residual dependence. By contrast, combinatorial algorithms with a criterion that is not score-equivalent (such as the MSE) depend on the data scale. Due to the optimization procedure, continuous structure learning algorithms may depend on the data scale irrespective of whether the employed score is score-equivalent (as, for example, $GOLEM$ for Gaussian models) or not (as, for example, $GOLEM$ under likelihood misspecification or $NOTEARS$).

We first establish how varsortability affects the gradients of MSE-based score functions (which are akin to assuming equal noise variances in the Gaussian setting) and when initializing with the empty graph $0_{d \times d}$ (as is done in $NOTEARS$ and $GOLEM$). Full statements of objective functions and respective gradients are found in Appendix E. Since $\nabla \text{MSE}(W) \propto \mathbf{X}^\top (\mathbf{X} - \mathbf{X}W)$ we have that $\nabla \text{MSE}(0_{d \times d}) \propto \mathbf{X}^\top \mathbf{X}$ and $\nabla \mathcal{L}_{EV}(0_{d \times d}) \propto 1/\|\mathbf{x}\|_2^2 \mathbf{X}^\top \mathbf{X}$. The initial gradient step of both $NOTEARS$ and $GOLEM-EV$ is symmetric. We have $\nabla \text{MSE}(W) \propto [\mathbf{X}^\top (x_1 - \mathbf{X}w_1), \dots, \mathbf{X}^\top (x_d - \mathbf{X}w_d)]$ where the j^{th} column $\mathbf{X}^\top (x_j - \mathbf{X}w_j)$ reflects the vector of empirical covariances of the j^{th} residual vector $x_j - \mathbf{X}w_j$ with each x_i . Provided a small identical step size is used across all entries of W in the first step (as, for example, in $GOLEM-EV$), we empirically find the residual variance after the first gradient step to be larger in those components that have higher marginal variance (see Appendix E.3 for a heuristic argument). We observe that during the next optimization steps $\nabla \text{MSE}(W)$ tends to be larger magnitude for edges pointing in the direction of nodes with high-variance residuals (which tends to be those with high marginal variance) than for those pointing in the direction of nodes with low-variance residuals (which tends to be those with low marginal variance). Intuitively, when cycles are penalized, the insertion of edges pointing to nodes with high residuals is favored as a larger reduction in MSE may be achieved than by including the opposing edge. Given high varsortability, this corresponds to favoring edges in the causal direction. This way, the global information about the causal order in case of high varsortability is effectively exploited.

Once we allow for unequal noise variances as in $GOLEM-NV$, the marginal variances lead the gradients differently. Letting $\text{MSE}_j(w_j) = \frac{1}{n} \|x_j - \mathbf{X}w_j\|_2^2$, we have

$$\nabla \mathcal{L}_{NV}(W) \propto \left[\frac{\mathbf{X}^\top (x_j - \mathbf{X}w_j)}{\text{MSE}_j(w_j)} \right]_{j=1, \dots, d}$$

such that the logarithmic derivative breaks the symmetry of the first step for the non-equal variance formulation of $GOLEM$ and we have $\nabla \mathcal{L}_{NV}(0_{d \times d}) \propto \mathbf{X}^\top \mathbf{X} \text{diag}(\|x_1\|_2^2, \dots, \|x_d\|_2^2)$. While $\nabla_W \text{MSE}(W) \propto [\mathbf{X}^\top (x_j - \mathbf{X}w_j)]_{j=1, \dots, d}$ tends to favor edges in causal direction (see above), the column-wise inverse MSE scaling of $\mathbf{X}^\top (x_j - \mathbf{X}w_j)$ by $\text{MSE}_j(w_j)$ (the residual variance in

the j^{th} component) leads to larger-magnitude gradient steps for edges pointing in the direction of low-variance nodes rather than high-variance nodes. Given high varsortability, this corresponds predominantly to the anti-causal direction.

We conjecture that the first gradient steps have a dominant role in determining the causal structure, even though afterwards the optimization is governed by a non-trivial interplay of optimizer, model fit, constraints, and penalties. For this reason we focus on the first optimization steps to explain a) why continuous structure learning algorithms that assume equal noise variance work remarkably well in the presence of high varsortability and b) why performance changes once data is standardized and the marginal variances no longer hold information about the causal order. Because of the acyclicity constraint, it may be enough for a weight $w_{i \rightarrow j}$ to be greater in magnitude than its counterpart $w_{j \rightarrow i}$ early on in the optimization for the smaller edge to be pruned from there on. For a discussion of the interplay between sparsity penalties and data scale see Appendix J, which indicates that the nodes need to be on a comparable data scale for l^1 -penalization to be well calibrated. Ng et al. [2020] provide further discussion on sparsity and acyclicity constraints in continuous DAG learning.

3.5 *sortnregress*: A Diagnostic Tool to Reveal Varsortability

We propose an algorithm *sortnregress* performing the following two steps:

order search Sort nodes by increasing marginal variance.

parent search Regress each node on all of its predecessors in that order, using a sparse regression technique to prune edges [Shojaie and Michailidis, 2010]. We employ Lasso regression [Tibshirani, 1996] using the Bayesian Information Criterion [Schwarz, 1978] for model selection.

As a baseline, *sortnregress* is easy to implement (cf. Appendix H.1) and highlights and evaluates to which extent the data scale is informative of the causal structure in different benchmark scenarios. An extension for non-linear additive noise models is obtained by using an appropriate non-linear regression technique in the parent search step, possibly paired with cross-validated recursive feature elimination. It facilitates a clear and contextualized assessment of different structure learning algorithms in different benchmark scenarios. The relationship between varsortability and the performance of *sortnregress* in a linear setting is shown in Appendix H.2. Varying degrees of varsortability and performance of *sortnregress* add an important dimension which current benchmarks do not consider.

4 Simulations

We compare the performance of the algorithms introduced in Section 2.3 on raw and standardized synthetic data. In our comparison, we distinguish between settings with different noise distributions, graph types, and graph sizes. Our experimental set-up follows those in Zheng et al. [2018], Ng et al. [2020] and we contribute results obtained repeating their experiments in Appendix K. We complement our and previous DAG-recovery results by additionally evaluating how well the DAG output by continuous structure learning algorithms identifies the MEC of the ground-truth DAG.

4.1 Data Generation

We sample Erdős-Rényi (ER) [Erdős and Rényi, 1960] and Scale-Free (SF) [Barabási and Albert, 1999] graphs and the parameters for ANMs according to the simulation details in Table 1. For a graph specified as ER- k or SF- k with d nodes, we simulate dk edges. For every combination of parameters, we create a raw data instance and a standardized version that is de-meaned and re-scaled to unit variance. On standardized data, we have varsortability $v = \frac{1}{2}$ and the marginal variances hold no information about the causal ordering of the nodes. In all our experimental settings, varsortability averages above 0.94 on the raw data scale (cf. Appendix G.1).

Table 1: Parameters for synthetic data generation.

Repetitions	10	Edge weights	iid Unif($(-2, -.5) \cup (.5, 2)$)
Graphs	ER-2, SF-2, SF-4	Noise distributions	Exponential, Gaussian, Gumbel
Nodes	$d \in \{10, 30, 50\}$	Noise standard deviations	1 (Gaussian-EV); iid Unif(.5, 2) (others)
Samples	$n = 1000$		

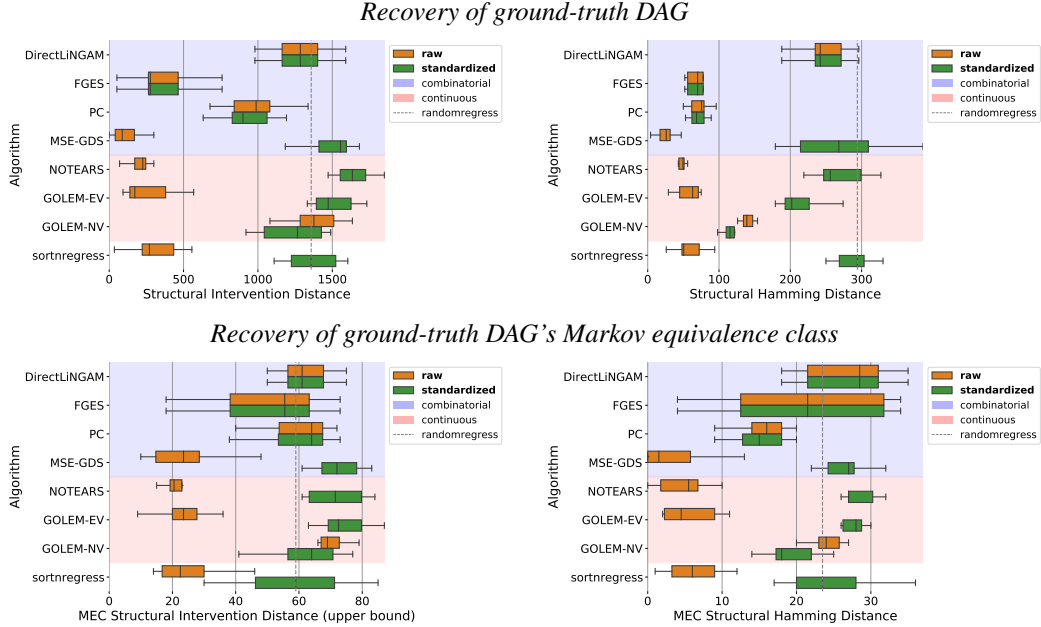


Figure 1: SID (left, lower is better) and SHD (right, lower is better) between recovered and ground-truth graphs (top) or Markov equivalence classes (bottom) for ER-2 graphs with 50 (top) or 10 (bottom) nodes and Gaussian-NV noise. The performance of *sortnregress*, which only exploits varsortability, matches that of the continuous methods *NOTEARS* and *GOLEM*.

4.2 Evaluation

We evaluate performance using structural intervention distance (SID) [Peters and Bühlmann, 2015] and structural Hamming distance (SHD) between the recovered graph and ground truth. Additionally, we contribute a performance assessments of continuous structure learning algorithms in terms of the SID and SHD between the ground-truth MEC and the recovered MEC (PC, FGES) or the MEC identified by the recovered graph (NOTEARS, GOLEM). SID assesses in how far the recovered structure enables to correctly predict the effect of interventions. SHD measures the distance between the true and recovered graph by counting how many edge insertions, deletions, and reversals are required to turn the former into the latter. Since interventional distributions can consistently be estimated given only the causal ordering [Bühlmann et al., 2014], SID is less susceptible to arbitrary choices of edge pruning procedures and thresholds than SHD. Intuitively, SID prioritizes the causal order, while SHD prioritizes the correctness of individual edges. We follow common practices for edge thresholding and scoring (see Appendix K).

4.3 Performance on Raw Versus Standardized Data

We group our algorithms into combinatorial and continuous algorithms. We propose a novel baseline algorithm termed *sortnregress* which serves as a reference marking the performance achievable by directly exploiting marginal variances. We indicate its performance on standardized data as *randomregress*, since it amounts to choosing a random order and regressing each node on its predecessors. We use boxplots aggregating the performance achieved in the 10 repetitions on raw and standardized data by each algorithm and create separate plots per noise type to account for identifiability and MLE specification differences. We show the results obtained for ER-2 graphs and Gaussian noise with non-equal variances in Figure 1. These results are representative of the results obtained for different graphs and noise distributions (cf. Appendix K). For the simulated settings, varsortability is high (> 0.94) on the raw data scale (cf. Appendix G.1).

We observe that some algorithms are highly scale-sensitive and perform vastly different on raw and standardized data. The algorithms *NOTEARS*, *MSE-GDS*, *GOLEM-EV* are most affected – their performance is excellent on raw data but far worse on standardized data. Note that all of these rely on

a loss function that revolves around the MSE. The performance of *GOLEM-NV* is also scale-sensitive but improves upon standardization. The direction of the effect of standardization is in line with the predictions by our gradient analysis in Section 3.4. Note that we initialize all algorithms with the empty graph since we are primarily interested in comparing the impact of standardization given equal starting conditions. On standardized data, an initialization of *GOLEM-NV* with the results of *GOLEM-EV*, as recommended by Ng et al. [2020], does not improve performance and may fail to converge. *sortnregress* achieves competitive performance on raw, and baseline performance on standardized data. It thus qualifies as diagnostic tool to highlight how much of a given causal structure learning task can be resolved by exploiting data scale and sorting nodes by their marginal variance.

In summary, the evidence corroborates our claim that the remarkable performance on raw data and the overall behavior upon standardization of the continuous structure learning algorithms may be driven primarily by high varsortability. On a real-world causal protein signaling dataset [Sachs et al., 2005] we measure a mean varsortability of 0.57 (which is close to chance level at 0.5) with a standard deviation of 0.01 across our bootstrapped samples and do not observe the consistent performance pattern described for synthetic data with high varsortability (cf. Appendix I).

5 Gaming Further Benchmarks

5.1 Orienting Causal Chains on Standardized Data

In order to design a causal discovery benchmark that does not favor methods that explicitly exploit marginal variances we may standardize the data or employ coefficient re-scaling schemes. Mooij et al. [2020], for example, propose a scale-harmonization by dividing each column $w_j = [w_{k \rightarrow j}]_{j=1, \dots, d} \in \mathbb{R}^d$ of the drawn adjacency matrices by $\sqrt{\|w_j\|^2 + 1}$ such that each variable would have comparable scale if all its direct parents were independently standard-normally distributed. However, this does not avoid the problem of potentially inadvertent patterns in simulated ANMs. Even after standardization or scale-harmonization, DAGs with previously high varsortability generate data with distinct covariance patterns that may be exploited.

In Appendix F we present an instructive example of a decision rule that can infer the orientation of a causal chain from raw, standardized, and scale-harmonized data with accuracy strictly greater than 50%. For a causal chain $X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_d$ where edge weights and noise terms are drawn iid we can decide between the two Markov-equivalent graphs $X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_d$ and $X_1 \leftarrow X_2 \leftarrow \dots \leftarrow X_d$ with above-chance accuracy. The empirical results for varying chain-lengths and various edge-weight distribution are deferred to the appendix where we discuss the 3-variable chain in detail and illustrate that the phenomenon extends from finite-sample to the population setting.

The intuition is as follows. Consider data generated by $X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_d$ and the aim is to infer from standardized data whether $X_1 \rightarrow \dots \rightarrow X_d$ or $X_1 \leftarrow \dots \leftarrow X_d$. For data with high varsortability and comparable noise variance on the raw data scale it holds that the further downstream a node X_i is in the causal chain, the stronger the variance of its parent $\text{Var}(X_{i-1})$ contributes to its marginal variance $\text{Var}(X_i) = \text{Var}(X_{i-1}) + \text{Var}(N_i)$ relative to its noise variance $\text{Var}(N_i)$, and the stronger is it correlated with its parent. Thus, the sequence of regression coefficients, which in the standardized case amounts to $(\text{Corr}(X_i, X_{i+1}))_{i=1, \dots, d-1}$, tends to increase in magnitude along the causal order and decrease in the anti-causal direction. The proposed decision rule predicts the causal direction as the one in which the absolute values of the regression coefficients tend to increase. This chain orientation rule achieves above-chance performance on raw, standardized, and scale-harmonized data (cf. Appendix F).

5.2 Sorting by Variance in Non-Linear Settings

Varsortability may also be exploited in non-linear settings. Table 2 shows the results of sorting by marginal variance and filling in all edges from lower-variance nodes to higher-variance nodes in a non-linear setting. This variance sorting strategy is more naive than *sortnregress* and places no assumption on the functional form. The results are substantially better than random sorting and may therefore be a more informative baseline than commonly used random graphs. We do not show performance in terms of SHD, as our variance sorting baseline always yields a fully connected graph. Although the data generating process is not identical, we note that the improvement of our crude variance sorting over random sorting compares favorably to some of the improvements gained by

more involved methods over random graphs as shown in [Lachapelle et al. \[2019, Table 1\]](#). Our results indicate that exploiting varsortability may also deliver competitive results in non-linear settings.

Table 2: SID of naive baselines on non-linear data. Results on 1000 observations of additive Gaussian process ANMs with noise variance 1 simulated as by [Zheng et al. \[2020\]](#) (10 repetitions each; average varsortability ν per graph type shown in parentheses).

Algorithm	Graph (average ν)	ER-1 (0.87)	ER-4 (0.95)	SF-1 (0.95)	SF-4 (0.98)
variance sorting		7.7 ± 5.72	25.2 ± 12.36	1.9 ± 2.28	7.6 ± 3.37
random sorting		27.9 ± 11.44	63.1 ± 8.10	22.3 ± 13.14	59.5 ± 7.32

We find similarly high levels of varsortability for many non-linear functional relationships and graph parameters (cf. Appendix G.2). This begs the question how much other successful methods exploit varsortability, how they compare to non-linear nonparametric methods that leverage assumptions on the residual variances [[Gao et al., 2020](#)], and how they perform under data standardization. We encourage such an exploration in future work and suggest that varsortability and *sortnregress* or *variance sorting* should always be included in future benchmarks.

6 Discussion and Conclusion

We find that continuous structure learning methods are highly susceptible to data rescaling and some do not perform well without access to the true data scale. Therefore, scale-variant causal structure learning methods should be applied and benchmarked with caution, especially if the variables do not share a measurement scale or when the true scale of the data is unattainable. It is important to declare whether data is standardized prior to being fed to various structure learning algorithms.

Following the first release of the present paper, [Kaiser and Sipos \[2021\]](#) also independently reported the drop in performance of *NOTEARS* upon standardizing the data and presented a low-dimensional exemplary case. Beyond a reporting of impaired *NOTEARS* performance, we also analyze score-equivalent methods, provide exhaustive simulation experiments, and explain the phenomenon.

Our aim is to raise awareness of the severity with which scaling properties in data from simulated DAGs and causal additive models may distort algorithm performance. Increasing marginal variances can render scenarios identifiable, which may commonly not be expected to be so—for example the Gaussian case with non-equal variances. We therefore argue that varsortability should be taken into account for future benchmarking. Yet, with any synthetic benchmark there remains a risk that the results are not indicative of algorithm performance on real data. Our results indicate that current structure learning algorithms may perform within the range of naive baselines on real-world datasets.

The theoretical results of our paper are limited to the setting of linear ANMs. Additionally, our conjecture regarding the importance of the first gradient steps, and with it a rigorous causal explanation for the learning behavior of different continuous algorithms and corresponding score functions remain open and require further research to be settled. Our empirical findings indicate that causal discovery benchmarks can be similarly gamed on standardized data and in non-linear settings, but further research is needed to confirm this. We focus on a specific subset of algorithms, the impact of patterns in benchmarking data on a wider class of algorithms and score functions remains to be explored.

Varsortability arises in many ANMs and the marginal variances increase drastically along the causal order, at least in common simulation settings. This begs the question what degree of varsortability can be observed or assumed in real-world data. If the marginal variances carry information about the causal order, our results suggest that it can and should be leveraged for structure learning. Otherwise, our contribution motivates future research into representative benchmarks and may put the practical applicability of the additive noise assumption into question.

Acknowledgements

We thank Jonas M. Kübler, Jonas Peters, and Sorawit Saengkyongam for helpful discussions and comments. SW was supported by the Carlsberg Foundation.

References

- Bryon Aragam and Qing Zhou. Concave penalized estimation of sparse Gaussian Bayesian networks. *Journal of Machine Learning Research*, 16(69):2273–2328, 2015.
- Albert-László Barabási and Réka Albert. Emergence of scaling in random networks. *Science*, 286(5439):509–512, 1999.
- Rohit Bhattacharya, Tushar Nagarajan, Daniel Malinsky, and Ilya Shpitser. Differentiable causal discovery under unmeasured confounding. In *International Conference on Artificial Intelligence and Statistics*, 2021.
- Philippe Brouillard, Sébastien Lachapelle, Alexandre Lacoste, Simon Lacoste-Julien, and Alexandre Drouin. Differentiable causal discovery from interventional data. In *Advances in Neural Information Processing Systems*, 2020.
- Peter Bühlmann, Jonas Peters, and Jan Ernest. CAM: Causal additive models, high-dimensional order search and penalized regression. *The Annals of Statistics*, 42(6):2526–2556, 2014.
- Wenyu Chen, Mathias Drton, and Y Samuel Wang. On causal discovery with an equal-variance assumption. *Biometrika*, 106(4):973–980, 09 2019.
- David M. Chickering. A Transformational Characterization of Equivalent Bayesian Network Structures. In *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence (UAI)*, 1995.
- David M. Chickering. Learning Bayesian Networks is NP-complete. In *Learning From Data*, pages 121–130. Springer, 1996.
- David M. Chickering. Learning equivalence classes of Bayesian-network structures. *Journal of Machine Learning Research*, 2:445–498, 2002a.
- David M. Chickering. Optimal structure identification with greedy search. *Journal of Machine Learning Research*, 3(11):507–554, 2002b.
- David M. Chickering, David Heckerman, and Chris Meek. Large-sample learning of Bayesian networks is NP-hard. *Journal of Machine Learning Research*, 5:1287–1330, 2004.
- A. Philip Dawid. Beware of the DAG! In *Proceedings of Workshop on Causality: Objectives and Assessment at NIPS 2008*, volume 6 of *Proceedings of Machine Learning Research*, pages 59–86. PMLR, 2010.
- Paul Erdős and Alfréd Rényi. On the evolution of random graphs. *Publication of the Mathematical Institute of the Hungarian Academy of Sciences*, 5(1):17–60, 1960.
- Ming Gao, Yi Ding, and Bryon Aragam. A polynomial-time algorithm for learning nonparametric causal graphs. In *Advances in Neural Information Processing Systems*, 2020.
- Asish Ghoshal and Jean Honorio. Learning Identifiable Gaussian Bayesian Networks in Polynomial Time and Sample Complexity. In *Advances in Neural Information Processing Systems*, 2017.
- Asish Ghoshal and Jean Honorio. Learning linear structural equation models in polynomial time and sample complexity. In *International Conference on Artificial Intelligence and Statistics*, 2018.
- Patrik Hoyer, Dominik Janzing, Joris M. Mooij, Jonas Peters, and Bernhard Schölkopf. Nonlinear causal discovery with additive noise models. In *Advances in Neural Information Processing Systems*, 2009.
- Guido W. Imbens and Donald B. Rubin. *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press, 2015.
- Marcus Kaiser and Maksim Sipos. Unsuitability of NOTEARS for Causal Graph Discovery. *arXiv preprint arXiv:2104.05441*, 2021.
- Sébastien Lachapelle, Philippe Brouillard, Tristan Deleu, and Simon Lacoste-Julien. Gradient-based neural DAG learning. *arXiv preprint arXiv:1906.02226*, 2019.
- Hao-Chih Lee, Matteo Danieletto, Riccardo Miotto, Sarah T. Cherng, and Joel T. Dudley. Scaling structural learning with NO-BEARS to infer causal transcriptome networks. *arXiv preprint arXiv:1911.00081*, 2019.
- Po-Ling Loh and Peter Bühlmann. High-dimensional learning of linear causal networks via inverse covariance estimation. *Journal of Machine Learning Research*, 15:3065–3105, 2014.

- Christopher Meek. Strong Completeness And Faithfulness In Bayesian Networks. In *Proceedings of Eleventh Conference on Uncertainty in Artificial Intelligence (UAI)*, 1995.
- Christopher Meek. *Graphical Models: Selecting causal and statistical models*. PhD thesis, Carnegie Mellon University, 1997.
- Joris M. Mooij, Sara Magliacane, and Tom Claassen. Joint Causal Inference from Multiple Contexts. *Journal of Machine Learning Research*, 21(99):1–108, 2020.
- Ignavier Ng, AmirEmad Ghassami, and Kun Zhang. On the role of sparsity and DAG constraints for learning linear DAGs. In *Advances in Neural Information Processing Systems*, 2020.
- Roxana Pamfil, Nisara Sriwattanaworachai, Shaan Desai, Philip Pilgerstorfer, Konstantinos Georgatzis, Paul Beaumont, and Bryon Aragam. DYNOTEARS: Structure Learning from Time-Series Data. In *International Conference on Artificial Intelligence and Statistics*, 2020.
- Gunwoong Park. Identifiability of additive noise models using conditional variances. *Journal of Machine Learning Research*, 21(75):1–34, 2020.
- Judea Pearl. Causal inference in statistics: An overview. *Statistics surveys*, 3:96–146, 2009.
- Jonas Peters and Peter Bühlmann. Identifiability of Gaussian structural equation models with equal error variances. *Biometrika*, 101(1):219–228, 2014.
- Jonas Peters and Peter Bühlmann. Structural intervention distance for evaluating causal graphs. *Neural Computation*, 27(3):771–799, 2015.
- Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of Causal Inference*. MIT Press, 2017.
- Joseph D. Ramsey, Kun Zhang, Madelyn Glymour, Ruben S. Romero, Biwei Huang, Imme Ebert-Uphoff, Savini Samarasinghe, Elizabeth A. Barnes, and Clark Glymour. TETRAD—A toolbox for causal discovery. In *8th International Workshop on Climate Informatics (CI 2018)*, 2018.
- Kenneth J. Rothman, Sander Greenland, and Timothy L. Lash. *Modern Epidemiology*. Lippincott Williams & Wilkins, 2008.
- Jakob Runge, Xavier-Andoni Tibau, Matthias Bruhns, Jordi Muñoz-Marí, and Gustau Camps-Valls. The causality for climate competition. In *NeurIPS 2019 Competition and Demonstration Track*, pages 110–120. PMLR, 2020.
- Karen Sachs, Omar Perez, Dana Pe’er, Douglas A. Lauffenburger, and Garry P. Nolan. Causal protein-signaling networks derived from multiparameter single-cell data. *Science*, 308(5721):523–529, 2005.
- Andrew D. Sanford and Imad A. Moosa. A Bayesian network structure for operational risk modelling in structured finance operations. *Journal of the Operational Research Society*, 63(4):431–444, 2012.
- Bernhard Schölkopf. Causality for machine learning. *arXiv preprint arXiv:1911.10500*, 2019.
- Gideon Schwarz. Estimating the dimension of a model. *The Annals of Statistics*, 6(2):461–464, 1978.
- Shohei Shimizu, Patrik O. Hoyer, Aapo Hyvärinen, and Antti Kerminen. A linear non-Gaussian acyclic model for causal discovery. *Journal of Machine Learning Research*, 7:2003–2030, 2006.
- Shohei Shimizu, Takanori Inazumi, Yasuhiro Sogawa, Aapo Hyvärinen, Yoshinobu Kawahara, Takashi Washio, Patrik O. Hoyer, and Kenneth Bollen. DirectLiNGAM: A direct method for learning a linear non-Gaussian structural equation model. *Journal of Machine Learning Research*, 12:1225–1248, 2011.
- Ali Shojaie and George Michailidis. Penalized likelihood methods for estimation of sparse high-dimensional directed acyclic graphs. *Biometrika*, 97(3):519–538, 07 2010.
- Peter Spirtes and Clark Glymour. An algorithm for fast recovery of sparse causal graphs. *Social Science Computer Review*, 9(1):62–72, 1991.
- Peter Spirtes, Clark N. Glymour, Richard Scheines, and David Heckerman. *Causation, Prediction, and Search*. MIT press, 2000.
- Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.

- Matthew J. Vowels, Necati Cihan Camgoz, and Richard Bowden. D’ya like DAGs? A Survey on Structure Learning and Causal Discovery. *arXiv preprint arXiv:2103.02582*, 2021.
- Dennis Wei, Tian Gao, and Yue Yu. DAGs with No Fears: A Closer Look at Continuous Optimization for Learning Bayesian Networks. *arXiv preprint arXiv:2010.09133*, 2020.
- Sebastian Weichwald, Martin E. Jakobsen, Phillip B. Mogenssen, Lasse Petersen, Nikolaj Thams, and Gherardo Varando. Causal structure learning from time series: Large regression coefficients may predict causal links better in practice than small p-values. In *NeurIPS 2019 Competition and Demonstration Track*, pages 27–36. PMLR, 2020.
- Shulin Yang and Kuo-Chu Chang. Comparison of score metrics for Bayesian network learning. *IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and Humans*, 32(3):419–428, 2002.
- Yue Yu, Jie Chen, Tian Gao, and Mo Yu. DAG-GNN: DAG structure learning with graph neural networks. *arXiv preprint arXiv:1904.10098*, 2019.
- Jiji Zhang and Peter Spirtes. Detection of unfaithfulness and robust causal inference. *Minds and Machines*, 18(2):239–271, 2008.
- Kun Zhang and Aapo Hyvarinen. On the identifiability of the post-nonlinear causal model. In *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence (UAI)*, 2009.
- Xun Zheng, Bryon Aragam, Pradeep K. Ravikumar, and Eric P. Xing. DAGs with NO TEARS: Continuous optimization for structure learning. In *Advances in Neural Information Processing Systems*, 2018.
- Xun Zheng, Chen Dan, Bryon Aragam, Pradeep K. Ravikumar, and Eric P. Xing. Learning sparse nonparametric DAGs. In *International Conference on Artificial Intelligence and Statistics*, 2020.

Appendix

Table of Contents

A Varsortability in the Two-Node Case	15
B Derivation of Lower Bound on Pairwise Varsortability	15
C Varsortability and Identifiability by Conditional Noise Variances	16
C.1 Park, 2020, Theorem 4 conditions satisfied without varsortability	16
C.2 Varsortability without Park, 2020, Theorem 4 conditions satisfied	16
D Algorithms	17
E The Subtle Interplay Between Marginal Variance and Gradient Directions	18
E.1 Example	18
E.2 Stepwise Gradient Derivation	20
E.3 Increasing Marginal and Residual Variances	21
E.4 Gradient Asymmetry	21
E.5 NOTEARS	22
F Standardization Is Not Enough and Regression Coefficients Tend to Increase Along the Causal Order	22
F.1 Infinite Sample	22
F.2 Finite Sample	24
G Empirical Evaluation of Varsortability	25
G.1 Varsortability in Linear Additive Noise Models	25
G.2 Varsortability in Non-Linear Additive Noise Models	26
G.3 Causal Order and Marginal Variance	26
G.4 Varsortability Algorithm	27
H <i>sortnregress</i>: A Diagnostic Tool to Reveal Varsortability	28
H.1 Implementation of Sortnregress	28
H.2 Varsortability and Score Attainable by Variance Ordering	28
I Evaluation on Real-World Data	29
J Model Selection in Continuous Optimization	29
K Detailed Results	30
K.1 MEC Recovery	30
K.2 Results Across Thresholding Regimes	31
K.3 Results Across Noise Distributions and Graph Types	31
K.4 Results Across Noise Distributions, Graph Types, and Graph Sizes	33

A Varsortability in the Two-Node Case

Consider the following ground truth and two competing linear acyclic models:

Ground-truth model:	Model $M1$) “ $A \rightarrow B$ ”:	Model $M2$) “ $A \leftarrow B$ ”:
$A := N_A$	$\hat{A} = 0$	$\hat{A} = \hat{v}B$
$B := wA + N_B$	$\hat{B} = \hat{w}A$	$\hat{B} = 0$

where $w \neq 0$, N_A and N_B are independent zero-centred noise terms that follow some distributions with non-vanishing corresponding variance V_A and V_B . The model parameters $\hat{w} = \frac{\text{Cov}(A,B)}{\text{Var}(A)} = w$ and $\hat{v} = \frac{\text{Cov}(A,B)}{\text{Var}(B)} = \frac{wV_A}{\text{Var}(B)}$ are the corresponding ordinary least-squares linear regression coefficients.

We evaluate in which cases the true model $M1$ obtains a smaller MSE than the wrong model $M2$, to decide if and under which conditions a MSE-based orientation rule recovers the ground-truth edge direction:

$$\begin{aligned}
& \text{MSE}(M1) < \text{MSE}(M2) \\
\iff & \text{Var}(A) + \text{Var}(B - \hat{w}A) < \text{Var}(A - \hat{v}B) + \text{Var}(B) \\
\iff & \text{Var}(A) + \text{Var}((wN_A + N_B) - wN_A) < \text{Var}\left(N_A - \frac{wV_A}{\text{Var}(B)}(wN_A + N_B)\right) + \text{Var}(B) \\
\iff & V_A + V_B < \frac{V_A V_B}{\text{Var}(B)} + w^2 V_A + V_B \\
\iff & 0 < V_A \left(\frac{V_B}{\text{Var}(B)} - 1\right) + w^2 V_A \\
\iff & 0 < \frac{-w^2 V_A}{\text{Var}(B)} + w^2 \\
\iff & V_A < \text{Var}(B) \\
\iff & (1 - w^2)V_A < V_B
\end{aligned}$$

For error variances $V_A \leq V_B$ and any non-zero edge weight w , the MSE-based inference is correct. This resembles known scale-based identifiability results based on equal or monotonically increasing error variances [Peters and Bühlmann, 2014, Park, 2020].

B Derivation of Lower Bound on Pairwise Varsortability

Let A and B be any two nodes in the sampled graph with edge $A \xrightarrow{w} B$, noise terms N_A, N_B , and without common ancestors and no other directed path from A to B . When sampling edge coefficients and noise variances randomly for the simulation of ANMs, distributions are incurred over the variances of A and B across those simulated ANMs. Let edge weights be sampled as $[W_{x \rightarrow y}]_{x,y=1,\dots,d} \sim \mathbb{P}_W$, and noise variances be sampled as $[S_{N_y}^2]_{y=1,\dots,d} \sim \mathbb{P}_{S^2}$. Across simulations, the marginal variances of A and B are transformations of S and W and themselves random variables denoted as V_A and V_B . The marginal variance V_Y of any node Y depends on its noise variance and the additional variance incurred by predecessor nodes given as $\sum_{X \in \text{PA}(Y)} W_{X \rightarrow Y}^2 V_X$. We can therefore bound the probability for the variable pair (A, B) to be varsortable from below via

$$\begin{aligned}
\text{P}[V_A < V_B] &= \text{P}\left[V_A < \left(W_{AB}^2 V_A + \sum_{X \in \text{PA}(B) \setminus \{A\}} W_{XB}^2 V_X + S_{N_B}^2\right)\right] \\
&\geq \text{P}[V_A < W_{AB}^2 V_A + S_{N_B}^2]
\end{aligned}$$

where equality holds if A is the only parent of B contributing to B 's marginal variance.

In common benchmarks, edge weights are drawn independently according to $\mathbb{P}_W \sim \otimes_{k,j=1,\dots,d} \text{Unif}((-2, -.5) \cup (.5, 2))$ and noise standard deviations are drawn iid $S_{N_j} \sim \text{Unif}(.5, 2)$.

C Varsortability and Identifiability by Conditional Noise Variances

While closely related, varsortability is not equivalent to the identifiability conditions laid out in Theorem 4, Park [2020], (henceforth referred to as ‘‘Theorem 4’’). We prove this by providing two examples. In Appendix C.1 part A) of the conditions in Theorem 4 is satisfied, while varsortability does not hold. In Appendix C.2 varsortability holds but neither part A) nor part B) of Theorem 4 are satisfied.

C.1 Park, 2020, Theorem 4 conditions satisfied without varsortability

Consider the following ground-truth model with unique causal order A, B, C :

$$\begin{aligned} A &:= N_A \\ B &:= \beta_{A \rightarrow B} A + N_B = 1A + N_B \\ C &:= \beta_{B \rightarrow C} B + N_C = \sqrt{\frac{2}{3}} B + N_C \end{aligned}$$

where N_A, N_B, N_C are jointly independent zero-centred noise terms with respective variances $\sigma_A^2 = 4, \sigma_B^2 = 2, \sigma_C^2 = 1$. The marginal variances are $\text{Var}(A) = 4 < \text{Var}(C) = 5 < \text{Var}(B) = 6$. Our example resembles the examples in Section 3.1 of Park [2020]. We can verify the three conditions for part A) of Theorem 4:

$$\begin{aligned} (A1) \quad & \sigma_A^2 < \sigma_B^2 + \beta_{A \rightarrow B}^2 \sigma_A^2, \\ (A2) \quad & \sigma_B^2 < \sigma_C^2 + \beta_{B \rightarrow C}^2 \sigma_B^2, \\ (A3) \quad & \sigma_A^2 < \sigma_C^2 + \beta_{B \rightarrow C}^2 \sigma_B^2 + \beta_{A \rightarrow B}^2 \beta_{B \rightarrow C}^2 \sigma_A^2 \end{aligned}$$

Inserting the values from above, we obtain

$$(A1) \quad 4 < 2 + 1 \cdot 4, \quad (A2) \quad 2 < 1 + \frac{2}{3} \cdot 2, \quad (A3) \quad 4 < 1 + \frac{2}{3} \cdot 2 + 1 \cdot \frac{2}{3} \cdot 4.$$

Our result verifies that identifiability is given as per Theorem 4 in Park [2020], while the order of increasing marginal variances is not in complete agreement with the causal order and varsortability is not equal to 1.

C.2 Varsortability without Park, 2020, Theorem 4 conditions satisfied

Consider the following ground-truth model with unique causal order A, B, C :

$$\begin{aligned} A &:= N_A \\ B &:= \beta_{A \rightarrow B} A + N_B = A + N_B \\ C &:= \beta_{A \rightarrow C} A + \beta_{B \rightarrow C} B + N_C = \frac{1}{\sqrt{2}} A + \frac{1}{\sqrt{2}} B + N_C \end{aligned}$$

where N_A, N_B, N_C are jointly independent zero-centred noise terms with respective variances $\sigma_A^2 = 4, \sigma_B^2 = 3, \sigma_C^2 = 1$. The marginal variances are $\text{Var}(A) = 4 < \text{Var}(B) = 7 < \text{Var}(C) = 10.5$. We now verify, that for both case A) and B) in Theorem 4 of Park [2020] at least one of the inequality constraints is violated.

One of the three conditions in A) is

$$\sigma_B^2 < \sigma_C^2 + \beta_{B \rightarrow C}^2 \sigma_B^2,$$

while for the above model we have

$$3 \not\prec 1 + \frac{1}{2} \cdot 3.$$

One of the three conditions in B) is

$$\frac{\sigma_C^2}{\sigma_B^2} > (1 - \beta_{B \rightarrow C}^2),$$

while for the above model we have

$$\frac{1}{3} \not\prec (1 - \frac{1}{2}).$$

For both criteria A) and B) in Theorem 4 at least one of the inequalities is not satisfied. We thus verify that even if identifiability is not given as per the sufficient conditions in Theorem 4, Park [2020], vartortability may still render the causal order identifiable.

D Algorithms

DirectLiNGAM is a method for learning linear non-Gaussian acyclic models [Shimizu et al., 2011]. It recovers the causal order by iteratively selecting the node whose residuals are least dependent on any predecessor node. In a strictly non-Gaussian setting, *DirectLiNGAM* is guaranteed to converge to the optimal solution asymptotically within a small fixed number of steps and returns a DAG. We use the implementation provided by the authors¹. We deliberately keep the default of a least-angle regression penalized by the Bayesian Information Criterion. We find that this penalty strikes a good balance between SID and SHD performance. Cross-validated least-angle regression performs better in terms of SID but poorer in terms of SHD.

PC [Spirtes and Glymour, 1991] is provably consistent in estimating the Markov equivalence class of the true data-generating graph if the causal Markov and faithfulness assumptions hold. The algorithm returns a completed partially directed acyclic graph (CPDAG). For computational reasons, we refrain from computing the lower and upper bounds of the SID for comparing CPDAGS with the ground-truth DAG as proposed by Peters and Bühlmann [2015]. Instead, we adopt the approach by Zheng et al. [2018] and resolve bidirectional edges favorably to obtain a DAG. We use the implementation in the *Tetrad*² package Ramsey et al. [2018].

FGES is an optimized version of the fast greedy equivalence search algorithm developed by Meek [1997], Chickering [2002b]. Under causal Markov and faithfulness assumptions, it is provably consistent for estimating the Markov equivalence class of the true data-generating graph. The algorithm returns a CPDAG, which we resolve favorably to obtain a DAG. We use the implementation in the *Tetrad*³ package [Ramsey et al., 2018].

MSE-GDS is a greedy DAG search procedure with a MSE score criterion. We implement *MSE-GDS* following other GDS procedures, for example, as described by Peters and Bühlmann [2014, Section 4], but use the MSE as score criterion instead of a likelihood- or BIC-based score criterion. For simplicity and computational ease, we consider a smaller search space and greedily forward-search over new edge insertions only instead of greedily searching over all neighbouring DAGs obtainable by edge insertions, removals, and deletions. For the linear setting, linear regression is used to determine the edge weights and the corresponding MSE-score for a given graph. For the non-linear setting, support vector regression can be used instead. The algorithm returns a DAG.

NOTEARS is a score-based method that finds both structure and parameters simultaneously by continuous optimization [Zheng et al., 2018]. The optimization formulation is based on the mean squared error and includes a sparsity penalty parameter λ and a differentiable acyclicity constraint:

$$\operatorname{argmin}_{W \in \mathbb{R}^{d \times d}} \operatorname{MSE}_{\mathbf{X}}(W) + \lambda \|W\|_1 \quad \text{s.t.} \quad \operatorname{tr}(\exp(W \odot W)) - d = 0.$$

¹<https://github.com/cdt15/lingam>

²<https://github.com/cmu-phil/tetrad>

³<https://github.com/cmu-phil/tetrad>

The algorithm returns a DAG. We use the implementation provided by the authors⁴. Throughout all our experiments we use *NOTEARS* over *NOTEARS-LI* (setting $\lambda = 0$), following the findings of Zheng et al. [2018, Tables 1 and 2, Figure 3], which suggest regularization only for samples smaller than the $n = 1000$ we use throughout.

GOLEM combines a soft version of the differentiable acyclicity constraint from Zheng et al. [2018] with a MLE objective [Ng et al., 2020]. The authors propose a multivariate Gaussian MLE for equal (EV) or unequal (NV) noise variances and optimize

$$\operatorname{argmin}_{W \in \mathbb{R}^{d \times d}} \tilde{\mathcal{L}}(W, \mathbf{X}) - \log(|\det(I - W)|) + \lambda_1 \|W\|_1 + \lambda_2 (\operatorname{tr}(\exp(W \odot W)) - d)$$

where $\tilde{\mathcal{L}}$ is either

$$\begin{aligned} \tilde{\mathcal{L}}_{EV}(W, \mathbf{X}) &= \frac{d}{2} (\mathcal{L}_{EV}(W, \mathbf{X}) + \log(n)) = \frac{d}{2} \log(n \operatorname{MSE}_{\mathbf{X}}(W)), \text{ or} \\ \tilde{\mathcal{L}}_{NV}(W, \mathbf{X}) &= \frac{1}{2} (\mathcal{L}_{NV}(W, \mathbf{X}) + d \log(n)) = \frac{1}{2} \sum_{j=1}^d \log(n \operatorname{MSE}_j(w_j)). \end{aligned}$$

We use the implementation and hyperparameters provided by the authors⁵. We train for 10^4 episodes as we found that half of that suffices to ensure convergence. Notably, we do not perform pretraining for our version of *GOLEM-NV*.

sortnregress is implemented as shown in Appendix H.1. We find that a least-angle regression penalized by the Bayesian Information Criterion strikes a good balance between SID and SHD performance.

E The Subtle Interplay Between Marginal Variance and Gradient Directions

We describe observations about the gradients involved in the optimization procedures of *NOTEARS* and *GOLEM-EV-NV*. We present an instructive example in Appendix E.1 and provide some intuition about how the adjacency matrix changes throughout the optimization. For convenience and reference we provide gradients of the individual terms involved in the respective objective functions (cf. Appendix E.2). In Appendix E.3 we argue why the nodes’ residual variances for the first gradient steps in an unconstrained optimization of MSE- or log-MSE-EV-based (*GOLEM-EV*) objective functions with acyclicity penalties tend to follow the same ordering as the nodes’ marginal variances. We analyze gradient symmetry and asymmetry in *GOLEM-EV-NV*’s gradient descent optimization under varsortability in Appendix E.4. While the intuition for small step size gradient-based unconstrained optimization partially carries over to the *NOTEARS* optimization procedure, here the interplay between varsortability and gradient directions is intricate due to a constrained optimization that is solved via the augmented Lagrangian method and dual descent with line-search instead of gradient descent as used in *GOLEM* [Zheng et al., 2018] (cf. Appendix E.5).

The heuristic arguments presented here are preliminary and aim to provide intuition. The optimization behaviour also heavily depends on the implementation of the optimization routine. For example, the original implementation of *NOTEARS* fixes the diagonal of W at zero and leverages curvature information (L-BFGS-B), while *GOLEM* updates all entries of W and employs learning rate optimizers. Future research is required to determine how precisely continuous structure learning algorithms achieve state-of-the-art results on highly varsortable data and, given our observations, we expect explanations to be specific to individual algorithms and their distinct implementations.

E.1 Example

The following example considers the population limit and illustrates a few intuitions about gradient based optimization and varsortability. Consider data is generated according to

$$\begin{pmatrix} X \\ Y \end{pmatrix} = \begin{pmatrix} 0 & \beta \\ 0 & 0 \end{pmatrix}^\top \begin{pmatrix} X \\ Y \end{pmatrix} + \begin{pmatrix} N_X \\ N_Y \end{pmatrix}$$

⁴<https://github.com/xunzheng/notears>

⁵<https://github.com/ignavier/golem>

where N_X and N_Y are independently normally distributed with standard deviations σ_{N_X} and σ_{N_Y} . Here, varsortability $v = 1$ and $1 = \text{Var } X < \text{Var } Y = 2$.

Initializing the weight matrix at the zero matrix, the gradient of the population MSE is

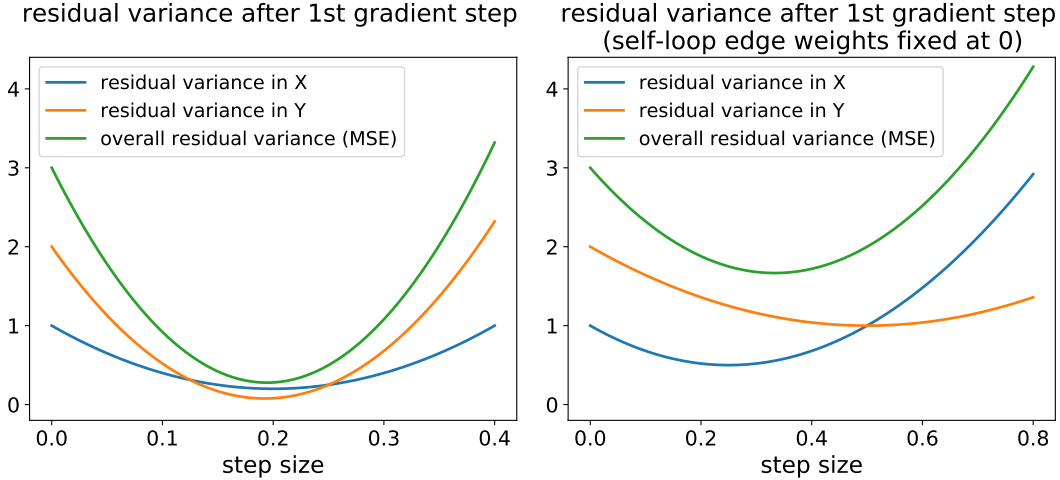
$$-2 \begin{pmatrix} \text{Var}(X) & \text{Cov}(X, Y) \\ \text{Cov}(Y, X) & \text{Var}(Y) \end{pmatrix} = -2 \begin{pmatrix} \sigma_{N_X}^2 & \beta \sigma_{N_X}^2 \\ \beta \sigma_{N_X}^2 & \beta^2 \sigma_{N_X}^2 + \sigma_{N_Y}^2 \end{pmatrix}$$

(see also Appendix E.2). The models for X and Y after a first gradient descent step of step size η are

$$\begin{aligned} \hat{X} &= 2\eta(\sigma_{N_X}^2 X + \beta \sigma_{N_X}^2 Y) \\ \hat{Y} &= 2\eta(\beta \sigma_{N_X}^2 X + (\beta^2 \sigma_{N_X}^2 + \sigma_{N_Y}^2) Y) \end{aligned}$$

If the diagonal of the weight matrix is clamped to 0 throughout the optimization, the terms corresponding to self-loops ($2\eta\sigma_{N_X}^2 X$ in \hat{X} and $(\beta^2 \sigma_{N_X}^2 + \sigma_{N_Y}^2) Y$ in \hat{Y}) are dropped above. This is the case in the original implementation of NOTEARS, where the unconstrained subproblem is optimized via L-BFGS-B with identity bounds on the diagonal entries of W .

Below we visualize $\text{Var}(X - \hat{X})$ (residual variance in X), $\text{Var}(Y - \hat{Y})$ (residual variance in Y), and the MSE $\text{Var}(X - \hat{X}) + \text{Var}(Y - \hat{Y})$, for varying step sizes η of the first gradient step where we exemplarily choose $\beta = \sigma_{N_X} = \sigma_{N_Y} = 1$.



Since the residual variances change continuously for increasing step sizes, the residual variances follow the order of the marginal variances for small step sizes (cf. also Appendix E.3). Since in GOLEM we solve an unconstrained optimization problem by gradient descent (with small step size and learning rate), the order of residual variances tends to remain unchanged during the first optimization steps. The order of the residual variances may swap relative to the order of marginal variances, though, if line-search is employed to determine the step size that minimizes the MSE-objective. This is the case in NOTEARS, where the MSE is minimized by a dual descent routine with increasing weight on the acyclicity penalty term. Here, the first symmetric update of the weight matrix occurs with a large step size that minimizes the MSE (minimum of the green curves in above plots). The ordering of the resulting residual variances is less obvious. In the above example, if the diagonal terms of the weight matrix are updated as well (left), the residual variance order after the first gradient step is opposite to the marginal variance order. If the diagonal entries are clamped at 0 (as is the case in NOTEARS and corresponding to the setting shown on the right), the first gradient step in the above example leads to a scenario where the residual variance order follows the marginal variance order and where the resulting edge weight for the direction $X \leftarrow Y$ overshoots the optimum, that is, the blue curve's minimum is attained for a smaller step size than the green curve's minimum. The intuition is as follows: If we minimize the MSE the step size calibrates a trade-off between residual variances in the different nodes; the high marginal variance nodes dominate the MSE such that the step size that minimizes the MSE may result in ill-chosen edge weights for the edges incoming into low-variance nodes. In the next optimization step, the gradient of the MSE loss for the edge $X \rightarrow Y$ pushes towards increasing that edge weight, while it pushes for decreasing the edge weight

$X \leftarrow Y$ (besides a gradient contribution from the acyclicity constraint). As a result, the edge weights for $X \rightarrow Y$ and $X \leftarrow Y$ are equal after the first step of *NOTEARS*, but better calibrated for the direction from low- to high-variance nodes, which here corresponds to the correct edge $X \rightarrow Y$. In the subsequent optimization step, decreasing the edge weight $X \leftarrow Y$ is favored both by the MSE gradient and the acyclicity penalty, while for the correct edge $X \rightarrow Y$ the MSE gradient pushes to further increasing the edge. Intuitively, if one needs to cut one of the two edges to avoid cycles, it is “cheaper” in MSE to cut the wrong edge $X \leftarrow Y$ from a high- to low-variance node.

E.2 Stepwise Gradient Derivation

MSE For $\mathbf{X} \in \mathbb{R}^{n \times d}$, the gradient of $\text{MSE}_{\mathbf{X}}(W) = \frac{1}{n} \|\mathbf{X} - \mathbf{X}W\|_2^2$ is

$$\begin{aligned} \nabla_W \text{MSE}_{\mathbf{X}}(W) &= \frac{1}{n} \nabla_W (\text{Tr}[\mathbf{X}^\top \mathbf{X}] - \text{Tr}[W^\top \mathbf{X}^\top \mathbf{X}] - \text{Tr}[\mathbf{X}^\top \mathbf{X}W] + \text{Tr}[W^\top \mathbf{X}^\top \mathbf{X}W]) \\ &= \frac{1}{n} (-\mathbf{X}^\top \mathbf{X} - \mathbf{X}^\top \mathbf{X} + \mathbf{X}^\top \mathbf{X}W + \mathbf{X}^\top \mathbf{X}W) \\ &= -\frac{2}{n} (\mathbf{X}^\top \mathbf{X} - \mathbf{X}^\top \mathbf{X}W) \\ &= -\frac{2}{n} \mathbf{X}^\top (\mathbf{X} - \mathbf{X}W) \\ &\propto \mathbf{X}^\top (\mathbf{X} - \mathbf{X}W) \end{aligned}$$

If W is polynomial in $\mathbf{X}^\top \mathbf{X}$, $\nabla_W \text{MSE}_{\mathbf{X}}(W)$ is symmetric. $\nabla_W \text{MSE}_{\mathbf{X}}(\mathbf{0}_{d \times d}) = -\frac{2}{n} \mathbf{X}^\top \mathbf{X}$.

GOLEM-EV The gradient of the unnormalized negative likelihood-part of the *GOLEM-EV* objective denoted as $\tilde{\mathcal{L}}_{EV}(W, \mathbf{X})$ is

$$\begin{aligned} \nabla_W \tilde{\mathcal{L}}_{EV}(W, \mathbf{X}) &= \frac{d}{2} \nabla_W \log(n \text{MSE}_{\mathbf{X}}(W)) \\ &= \frac{d}{2} \frac{1}{\text{MSE}_{\mathbf{X}}(W)} \nabla_W \text{MSE}_{\mathbf{X}}(W) \\ &\propto \frac{1}{\text{MSE}_{\mathbf{X}}(W)} \mathbf{X}^\top (\mathbf{X} - \mathbf{X}W) \end{aligned}$$

If W is polynomial in $\mathbf{X}^\top \mathbf{X}$, $\nabla_W \tilde{\mathcal{L}}_{EV}(W, \mathbf{X})$ is symmetric. $\nabla_W \tilde{\mathcal{L}}_{EV}(\mathbf{0}_{d \times d}, \mathbf{X}) = -\frac{d}{\|\mathbf{X}\|_2^2} \mathbf{X}^\top \mathbf{X}$.

GOLEM-NV The gradient of the unnormalized negative likelihood-part of the *GOLEM-NV* objective denoted as $\tilde{\mathcal{L}}_{NV}(W, \mathbf{X})$ is

$$\begin{aligned} \nabla_W \tilde{\mathcal{L}}_{NV}(W, \mathbf{X}) &= \frac{1}{2} \sum_{j=1}^d \nabla_W \log(n \text{MSE}_j(w_j)) \\ &= \left[-\frac{1}{n \text{MSE}_j(w_j)} \mathbf{X}^\top (x_j - \mathbf{X}w_j) \right]_{j=1, \dots, d} \\ &\propto \left[\frac{\mathbf{X}^\top (x_j - \mathbf{X}w_j)}{\text{MSE}_j(w_j)} \right]_{j=1, \dots, d} \end{aligned}$$

For the zero matrix, we have $\nabla_W \tilde{\mathcal{L}}_{NV}(\mathbf{0}_{d \times d}, \mathbf{X}) = -\mathbf{X}^\top \mathbf{X} \text{diag}(\|x_1\|_2^{-2}, \dots, \|x_d\|_2^{-2})$.

We focus on the gradients of MSE, \mathcal{L}_{EV} , and \mathcal{L}_{NV} since l1 penalty, acyclicity penalty h , LogDet term, and exact scaling of $\tilde{\mathcal{L}}_{EV}$ and $\tilde{\mathcal{L}}_{NV}$ play a subordinate role at the zero initialization, where the LogDet gradient has zero off-diagonals and $\nabla_W h$ vanishes:

The LogDet in *GOLEM-EV* and *GOLEM-NV* $\text{LogDet}(W) = \log(\det(I - W))$ has gradient

$$\nabla_W \text{LogDet}(W) = -(I - W)^{-\top}$$

and vanishes when W is the adjacency matrix of a DAG [Ng et al., 2020]. If W is symmetric, $\nabla_W \text{LogDet}(W)$ is symmetric. For the zero matrix, we have $\nabla_W \text{LogDet}(\mathbf{0}_{d \times d}) = -I$.

Acyclicity Penalty/Constraint The function $h(W) = \text{tr}(\exp(W \odot W)) - d$ has gradient $\nabla_W h(W) = \exp(W \odot W)^\top \odot 2W$. The $h(W)=0$ -level set characterizes adjacency matrices of DAGs [Zheng et al., 2018]. If W is symmetric, $\nabla_W h(W)$ is symmetric. For the zero matrix, we have $h(\mathbf{0}_{d \times d}) = 0$ and $\nabla_W h(\mathbf{0}_{d \times d}) = \mathbf{0}_{d \times d}$.

E.3 Increasing Marginal and Residual Variances

We observe a strong positive correlation between the ordering by marginal variance and the ordering by residual variance after the first gradient step when minimizing a MSE- or likelihood-based objective function via gradient descent with small step size (as in *GOLEM-EV/NV*). For small step sizes and learning rates, marginal variance order and residual variance order are perfectly aligned for the first few optimization steps. Here we argue for a MSE-based loss function why the residual variance follows the order of increasing marginal variance after the first optimisation step with sufficiently small step size. Future work may investigate subsequent optimisation steps and the non-MSE terms of the objective functions.

Consider the data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$. Without loss of generality, we assume the columns are zero-centred and ordered such that the sequence of diagonal entries $\text{diag}(\mathbf{X}^\top \mathbf{X})$ is weakly monotonically increasing. The diagonal entries $\text{diag}(\mathbf{X}^\top \mathbf{X})$ correspond to (n -times) the marginal variances at step 0. After the first gradient step with step size α in direction $-\nabla_W \text{MSE}_{\mathbf{X}}(\mathbf{0}_{d \times d}) = \frac{2}{n} \mathbf{X}^\top \mathbf{X}$ (cf. Appendix E.2) the vector of (n -times) the residual variances is

$$\begin{aligned} \mathbf{R} &= \text{diag}([\mathbf{X} - a\mathbf{X}\mathbf{X}^\top \mathbf{X}]^\top [\mathbf{X} - a\mathbf{X}\mathbf{X}^\top \mathbf{X}]) \\ &= \text{diag}(\mathbf{D}) - 2a \text{diag}(\mathbf{D}^2) + a^2 \text{diag}(\mathbf{D}^3) \end{aligned}$$

where $\mathbf{D} = \mathbf{X}^\top \mathbf{X}$ and $a = \frac{2}{n}\alpha$. For each coordinate i the residual variance \mathbf{R}_i is a continuous function in a (and α). For $a = 0$ and every $i \in [1, \dots, d-1]$ we have $\mathbf{R}_{i+1} - \mathbf{R}_i = \mathbf{D}_{i+1} - \mathbf{D}_i \geq 0$ with strict inequality if the variable pair $i, i+1$ is varsortable. Due to continuity, for any pair of variables with unequal marginal variances, there exists a sufficiently small step size to ensure that the resulting residual variances follow the same order as the marginal variances.

E.4 Gradient Asymmetry

We combine what we laid out in Appendix E.2.

The GOLEM-EV optimization problem is

$$\underset{W}{\text{argmin}} \tilde{\mathcal{L}}_{EV}(W, \mathbf{X}) - \text{LogDet}(W) + \lambda_1 \|W\|_1 + \lambda_2 h(W)$$

with the following gradient of the objective function

$$-\frac{d}{n \text{MSE}_{\mathbf{X}}(W)} \mathbf{X}^\top (\mathbf{X} - \mathbf{X}W) + (I - W)^{-\top} + \lambda_1 W \odot |W| + \lambda_2 \exp(W \odot W)^\top \odot 2W$$

which at zero reduces to

$$-\frac{d}{\|\mathbf{X}\|_2^2} \mathbf{X}^\top \mathbf{X} + I.$$

The GOLEM-NV optimization problem is

$$\underset{W}{\text{argmin}} \tilde{\mathcal{L}}_{NV}(W, \mathbf{X}) - \text{LogDet}(W) + \lambda_1 \|W\|_1 + \lambda_2 h(W)$$

with the following gradient of the objective function

$$\left[-\frac{1}{n \text{MSE}_j(w_j)} \mathbf{X}^\top (x_j - \mathbf{X}w_j) \right]_{j=1, \dots, d} + (I - W)^{-\top} + \lambda_1 W \odot |W| + \lambda_2 \exp(W \odot W)^\top \odot 2W$$

which at zero reduces to

$$-\mathbf{X}^\top \mathbf{X} \text{diag}(\|x_1\|_2^{-2}, \dots, \|x_d\|_2^{-2}) + I.$$

The gradient in *GOLEM-EV* is symmetric at $\mathbf{0}_{d \times d}$ at the first gradient descent step, but not in general for later steps. The gradient in *GOLEM-NV* is in general not symmetric and at $\mathbf{0}_{d \times d}$ (at the first

gradient descent step) the gradients for edges incoming into a node are inversely scaled by its marginal variance; consequently, for weights $w_{i \rightarrow j}$ and $w_{j \rightarrow i}$ of opposing edges the first gradient step is larger magnitude for the direction with lower-variance end-node and $w_{i \rightarrow j}$ is preferred over $w_{j \rightarrow i}$ if the variance of X_i is higher than that of X_j . Under high-varsortability, the first *GOLEM-NV* gradient step thus tends to favor edges in anti-causal direction over those in causal direction.

E.5 NOTEARS

The *NOTEARS* optimization problem is $\operatorname{argmin}_W \frac{1}{2} \operatorname{MSE}_{\mathbf{X}}(W)$ s.t. $h(W) = 0$ which is solved via the augmented Lagrangian method and dual descent [Zheng et al., 2018] (we omit the penalty term for the *NOTEARS-II* variant). In the original implementation, the algorithm is initialized at $\mathbf{0}_{d \times d}$ and the diagonal of W is not updated but fixed to zero (this amounts to dual projected descent, where the adjacency matrix is projected onto the matrices with zero diagonal at each step avoiding self-loops per fiat).

The augmented Lagrangian

$$\frac{1}{2} \operatorname{MSE}_{\mathbf{X}}(W) + \frac{\rho}{2} h(W)^2 + \alpha h(W)$$

has gradient

$$-\frac{1}{n} \mathbf{X}^\top (\mathbf{X} - \mathbf{X}W) + (\rho h + \alpha) (\exp(W \odot W)^\top \odot 2W)$$

which at zero reduces to

$$-\frac{1}{n} \mathbf{X}^\top \mathbf{X}$$

The step size of the first gradient step in direction $\propto \mathbf{X}^\top \mathbf{X}$ is optimized by line-search to minimize the overall MSE. As seen in the example in Appendix E.1, the residual variances may or may not follow the order of the marginal variances after this first step due to the step size being larger than the small step size that would ensure agreement between the orders (cf. Appendix E.3). Nonetheless, the step size optimized by line-search aims to optimize the overall MSE which tends to favor a better fit for edges incoming into nodes with high-marginal variance. As a result, the first gradient step results in edge weights that are better calibrated for edges incoming into high-marginal variance nodes than into low-marginal variance nodes. In subsequent steps of the dual ascent procedure with increasing acyclicity penalty, the reduction of overall MSE stands at odds with satisfying the DAG constraints; it is then more costly in terms of MSE to change the weights for edges into high-marginal nodes than into low-marginal nodes such that predominantly the edges into low-variance nodes tend to be removed to eventually satisfy the acyclicity constraint. Under high varsortability, this amounts to a preference for causal edges.

F Standardization Is Not Enough and Regression Coefficients Tend to Increase Along the Causal Order

Code to reproduce the calculations and results in this section is available at <https://github.com/Scriddie/Varsortability>.

F.1 Infinite Sample

Here, we first discuss the three-variable case to complement the intuition provided in the main text. Consider the following ground-truth linear additive acyclic models, where the second model corresponds to a standardization of the first, and the third model corresponds to a re-scaled version of the first following Mooij et al. [2020]:

Raw ground-truth model	Standardized model	Scale-harmonized model
$A := N_A$	$A_s := A/\sqrt{\operatorname{Var}(A)}$	$A_m := N_A$
$B := \beta_{A \rightarrow B} A + N_B$	$B_s := B/\sqrt{\operatorname{Var}(B)}$	$B_m := \frac{\beta_{A \rightarrow B}}{\sqrt{\beta_{A \rightarrow B}^2 + 1}} A_m + N_B$
$C := \beta_{B \rightarrow C} B + N_C$	$C_s := C/\sqrt{\operatorname{Var}(C)}$	$C_m := \frac{\beta_{B \rightarrow C}}{\sqrt{\beta_{B \rightarrow C}^2 + 1}} B_m + N_C$

where, following common benchmark sampling schemes, $N_A, N_B,$ and N_C are independent zero-centred noise terms that follow some distributions with non-vanishing standard deviations $\sigma_A, \sigma_B,$ and σ_C sampled independently from $\text{Unif}(.5, 2)$ and where $\beta_{A \rightarrow B}$ and $\beta_{B \rightarrow C}$ are independently drawn from $\text{Unif}((-2, -.5) \cup (.5, 2))$. For any two nodes X and Y , $\beta_{X \rightarrow Y}$ denotes an underlying model parameter, while $\widehat{\beta}_{X \rightarrow Y}$ denotes the ordinary least-squares linear regression coefficient when regressing Y onto X which is given as $\widehat{\beta}_{X \rightarrow Y} = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$.

Given observations from a variable triplet (X, Y, Z) , the *causal chain orientation task* is to infer whether the data generating causal chain is $X \rightarrow Y \rightarrow Z$, that is, $(X, Y, Z) = (A, B, C)$ or $Z \leftarrow Y \leftarrow X$, that is, $(Z, X, Y) = (A, B, C)$. While both graphs are Markov equivalent, we can identify the correct orientation of the causal chain, for all three considered scaling regimes, with accuracy strictly greater than 50% by applying the following procedure:

Chain orientation rule:

- If $|\widehat{\beta}_{X \rightarrow Y}| < |\widehat{\beta}_{Y \rightarrow Z}|$ and $|\widehat{\beta}_{Z \rightarrow Y}| > |\widehat{\beta}_{Y \rightarrow X}|$, conclude $(X, Y, Z) = (A, B, C)$.
We conclude that $X \rightarrow Y \rightarrow Z$, if the regression coefficients are increasing in magnitude when regressing pairwise from “left to right”.
- If $|\widehat{\beta}_{X \rightarrow Y}| > |\widehat{\beta}_{Y \rightarrow Z}|$ and $|\widehat{\beta}_{Z \rightarrow Y}| < |\widehat{\beta}_{Y \rightarrow X}|$, conclude $(X, Y, Z) = (C, B, A)$.
We conclude that $X \leftarrow Y \leftarrow Z$, if the regression coefficients are increasing in magnitude when regressing pairwise from “right to left”.
- Otherwise, flip a coin to decide the orientation of the underlying causal chain.

For each data scale regime, we can obtain the population regression coefficients and express those in terms of the sampled model coefficients $\beta_{A \rightarrow B}, \beta_{B \rightarrow C}, \sigma_A, \sigma_B, \sigma_C$:

- Raw ground-truth model
 - “left to right”: $\widehat{\beta}_{A \rightarrow B} = \beta_{A \rightarrow B}$ and $\widehat{\beta}_{B \rightarrow C} = \beta_{B \rightarrow C}$
 - “right to left”: $\widehat{\beta}_{C \rightarrow B} = \frac{\beta_{B \rightarrow C}(\beta_{A \rightarrow B}^2 \sigma_A^2 + \sigma_B^2)}{\beta_{A \rightarrow B}^2 \beta_{B \rightarrow C}^2 \sigma_A^2 + \beta_{B \rightarrow C}^2 \sigma_B^2 + \sigma_C^2}$ and $\widehat{\beta}_{B \rightarrow A} = \frac{\beta_{A \rightarrow B} \sigma_A^2}{\beta_{A \rightarrow B}^2 \sigma_A^2 + \sigma_B^2}$
- Standardized model
 - “left to right”:
 $\widehat{\beta}_{A_s \rightarrow B_s} = \frac{\beta_{A \rightarrow B} \sigma_A^2}{\sqrt{\beta_{A \rightarrow B}^2 \sigma_A^2 + \sigma_B^2} \sqrt{\sigma_A^2}}$ and $\widehat{\beta}_{B_s \rightarrow C_s} = \frac{\beta_{B \rightarrow C} \sqrt{\beta_{A \rightarrow B}^2 \sigma_A^2 + \sigma_B^2}}{\sqrt{\beta_{A \rightarrow B}^2 \beta_{B \rightarrow C}^2 \sigma_A^2 + \beta_{B \rightarrow C}^2 \sigma_B^2 + \sigma_C^2}}$
 - “right to left”:
 $\widehat{\beta}_{C_s \rightarrow B_s} = \frac{\beta_{B \rightarrow C} \sqrt{\beta_{A \rightarrow B}^2 \sigma_A^2 + \sigma_B^2}}{\sqrt{\beta_{A \rightarrow B}^2 \beta_{B \rightarrow C}^2 \sigma_A^2 + \beta_{B \rightarrow C}^2 \sigma_B^2 + \sigma_C^2}}$ and $\widehat{\beta}_{B_s \rightarrow A_s} = \frac{\beta_{A \rightarrow B} \sigma_A^2}{\sqrt{\beta_{A \rightarrow B}^2 \sigma_A^2 + \sigma_B^2} \sqrt{\sigma_A^2}}$
- Scale-harmonized model
 - Regression coefficients “from left to right”:
 $\widehat{\beta}_{A_m \rightarrow B_m} = \frac{\beta_{A \rightarrow B}}{\sqrt{\beta_{A \rightarrow B}^2 + 1}}$ and $\widehat{\beta}_{B_m \rightarrow C_m} = \frac{\beta_{B \rightarrow C}}{\sqrt{\beta_{B \rightarrow C}^2 + 1}}$
 - Regression coefficients “from right to left”:
 $\widehat{\beta}_{C_m \rightarrow B_m} = \frac{\beta_{B \rightarrow C} (\beta_{B \rightarrow C}^2 + 1)^{1.5} (\beta_{A \rightarrow B}^2 \sigma_A^2 + \sigma_B^2 (\beta_{A \rightarrow B}^2 + 1))}{\beta_{A \rightarrow B}^2 \beta_{B \rightarrow C}^2 \sigma_A^2 (\beta_{B \rightarrow C}^2 + 1) + \beta_{B \rightarrow C}^2 \sigma_B^2 (\beta_{A \rightarrow B}^2 + 1) (\beta_{B \rightarrow C}^2 + 1) + \sigma_C^2 (\beta_{A \rightarrow B}^2 + 1) (\beta_{B \rightarrow C}^2 + 1)^2}$
and $\widehat{\beta}_{B_m \rightarrow A_m} = \frac{\beta_{A \rightarrow B} \sigma_A^2 \sqrt{\beta_{A \rightarrow B}^2 + 1}}{\beta_{A \rightarrow B}^2 \sigma_A^2 + \sigma_B^2 (\beta_{A \rightarrow B}^2 + 1)}$

We obtain the following probabilities by Monte Carlo approximation, resampling the 5 model parameters 100,000 times:

Table 3: Chain orientation results in the population limit.

Weight distribution	Chain orientation rule cases	
Unif((-2, .5) ∪ (.5, 2))	$P \left[\hat{\beta}_{A \rightarrow B} < \hat{\beta}_{B \rightarrow C} \text{ and } \hat{\beta}_{C \rightarrow B} > \hat{\beta}_{B \rightarrow A} \right]$	29.376%
	$P \left[\hat{\beta}_{A \rightarrow B} > \hat{\beta}_{B \rightarrow C} \text{ and } \hat{\beta}_{C \rightarrow B} < \hat{\beta}_{B \rightarrow A} \right]$	5.486%
	$P \left[\text{“orientation rule correct on raw data”} \right]$	61.945%
	$P \left[\hat{\beta}_{A_s \rightarrow B_s} < \hat{\beta}_{B_s \rightarrow C_s} \text{ and } \hat{\beta}_{C_s \rightarrow B_s} > \hat{\beta}_{B_s \rightarrow A_s} \right]$	73.181%
	$P \left[\hat{\beta}_{A_s \rightarrow B_s} > \hat{\beta}_{B_s \rightarrow C_s} \text{ and } \hat{\beta}_{C_s \rightarrow B_s} < \hat{\beta}_{B_s \rightarrow A_s} \right]$	26.819%
	$P \left[\text{“orientation rule correct on standardized data”} \right]$	73.181%
	$P \left[\hat{\beta}_{A_m \rightarrow B_m} < \hat{\beta}_{B_m \rightarrow C_m} \text{ and } \hat{\beta}_{C_m \rightarrow B_m} > \hat{\beta}_{B_m \rightarrow A_m} \right]$	31.631%
	$P \left[\hat{\beta}_{A_m \rightarrow B_m} > \hat{\beta}_{B_m \rightarrow C_m} \text{ and } \hat{\beta}_{C_m \rightarrow B_m} < \hat{\beta}_{B_m \rightarrow A_m} \right]$	17.318%
	$P \left[\text{“orientation rule correct on scale-harmonized data”} \right]$	57.1565%
Unif((-0.9, -0.5) ∪ (0.5, 0.9))	$P \left[\hat{\beta}_{A \rightarrow B} < \hat{\beta}_{B \rightarrow C} \text{ and } \hat{\beta}_{C \rightarrow B} > \hat{\beta}_{B \rightarrow A} \right]$	31.033%
	$P \left[\hat{\beta}_{A \rightarrow B} > \hat{\beta}_{B \rightarrow C} \text{ and } \hat{\beta}_{C \rightarrow B} < \hat{\beta}_{B \rightarrow A} \right]$	18.124%
	$P \left[\text{“orientation rule correct on raw data”} \right]$	56.454%
	$P \left[\hat{\beta}_{A_s \rightarrow B_s} < \hat{\beta}_{B_s \rightarrow C_s} \text{ and } \hat{\beta}_{C_s \rightarrow B_s} > \hat{\beta}_{B_s \rightarrow A_s} \right]$	62.231%
	$P \left[\hat{\beta}_{A_s \rightarrow B_s} > \hat{\beta}_{B_s \rightarrow C_s} \text{ and } \hat{\beta}_{C_s \rightarrow B_s} < \hat{\beta}_{B_s \rightarrow A_s} \right]$	37.769%
	$P \left[\text{“orientation rule correct on standardized data”} \right]$	62.231%
	$P \left[\hat{\beta}_{A_m \rightarrow B_m} < \hat{\beta}_{B_m \rightarrow C_m} \text{ and } \hat{\beta}_{C_m \rightarrow B_m} > \hat{\beta}_{B_m \rightarrow A_m} \right]$	30.025%
	$P \left[\hat{\beta}_{A_m \rightarrow B_m} > \hat{\beta}_{B_m \rightarrow C_m} \text{ and } \hat{\beta}_{C_m \rightarrow B_m} < \hat{\beta}_{B_m \rightarrow A_m} \right]$	20.607%
	$P \left[\text{“orientation rule correct on scale-harmonized data”} \right]$	54.709%
Unif((-0.9, -0.1) ∪ (0.1, 0.9))	$P \left[\hat{\beta}_{A \rightarrow B} < \hat{\beta}_{B \rightarrow C} \text{ and } \hat{\beta}_{C \rightarrow B} > \hat{\beta}_{B \rightarrow A} \right]$	32.480%
	$P \left[\hat{\beta}_{A \rightarrow B} > \hat{\beta}_{B \rightarrow C} \text{ and } \hat{\beta}_{C \rightarrow B} < \hat{\beta}_{B \rightarrow A} \right]$	24.012%
	$P \left[\text{“orientation rule correct on raw data”} \right]$	54.234%
	$P \left[\hat{\beta}_{A_s \rightarrow B_s} < \hat{\beta}_{B_s \rightarrow C_s} \text{ and } \hat{\beta}_{C_s \rightarrow B_s} > \hat{\beta}_{B_s \rightarrow A_s} \right]$	55.790%
	$P \left[\hat{\beta}_{A_s \rightarrow B_s} > \hat{\beta}_{B_s \rightarrow C_s} \text{ and } \hat{\beta}_{C_s \rightarrow B_s} < \hat{\beta}_{B_s \rightarrow A_s} \right]$	44.210%
	$P \left[\text{“orientation rule correct on standardized data”} \right]$	55.790%
	$P \left[\hat{\beta}_{A_m \rightarrow B_m} < \hat{\beta}_{B_m \rightarrow C_m} \text{ and } \hat{\beta}_{C_m \rightarrow B_m} > \hat{\beta}_{B_m \rightarrow A_m} \right]$	31.867%
	$P \left[\hat{\beta}_{A_m \rightarrow B_m} > \hat{\beta}_{B_m \rightarrow C_m} \text{ and } \hat{\beta}_{C_m \rightarrow B_m} < \hat{\beta}_{B_m \rightarrow A_m} \right]$	25.136%
	$P \left[\text{“orientation rule correct on scale-harmonized data”} \right]$	53.3655%

We draw edge weights independently from the uniform distribution indicated in the first column of Table 3 and noise standard-deviations $\sigma_A, \sigma_B, \sigma_C$ are drawn independently from Unif(.5, 2) in all cases. A 99% confidence interval for the orientation accuracy under random guessing is (49.593%, 50.407%). The orientation rule achieves above chance accuracy in all regimes.

F.2 Finite Sample

Given observations from (X_1, \dots, X_d) generated by a linear ANM with either $X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_d$ or $X_d \rightarrow X_{d-1} \rightarrow \dots \rightarrow X_1$, we can decide the directionality by identifying the direction in which the absolute values of the regression coefficients tend to increase. More precisely, we compare the sequences of absolute regression coefficients

“left-to-right regression coefficients” $|\hat{\beta}_{X_1 \rightarrow X_2}|, \dots, |\hat{\beta}_{X_{d-1} \rightarrow X_d}|$

to

“right-to-left regression coefficients” $|\hat{\beta}_{X_d \rightarrow X_{d-1}}|, \dots, |\hat{\beta}_{X_2 \rightarrow X_1}|$.

We infer $X_1 \rightarrow \dots \rightarrow X_d$ if the former is in better agreement with an ascending sorting than the latter and infer $X_d \rightarrow \dots \rightarrow X_1$ otherwise.

In the main text, we discussed the case for standardized data where the regression coefficients for any two nodes X_i and X_j are given as $|\text{Corr}(X_i, X_j)|$. We expect the sequence of absolute regression coefficients to increase along the causal order because the correlation between consecutive nodes tends to be higher further downstream as parent nodes contribute more to a nodes marginal variance relative to its noise term.

On the raw data scale, the sequences of regression coefficients are

$$\begin{aligned} \text{“left-to-right”} & \frac{\sqrt{\text{Var}(X_2)}}{\sqrt{\text{Var}(X_1)}} |\text{Corr}(X_1, X_2)|, \dots, \frac{\sqrt{\text{Var}(X_d)}}{\sqrt{\text{Var}(X_{d-1})}} |\text{Corr}(X_{d-1}, X_d)| \quad \text{and} \\ \text{“right-to-left”} & \frac{\sqrt{\text{Var}(X_{d-1})}}{\sqrt{\text{Var}(X_d)}} |\text{Corr}(X_{d-1}, X_d)|, \dots, \frac{\sqrt{\text{Var}(X_1)}}{\sqrt{\text{Var}(X_2)}} |\text{Corr}(X_1, X_2)|. \end{aligned}$$

On both raw and standardized data, we find that the direction in which absolute regression coefficients tend to increase most corresponds to the causal direction in more than 50% of cases. To quantify “increasingness” of sequences of absolute regression coefficients we count the number of correctly ordered pairs of regression coefficients, that is, how often a regression coefficient is smaller in magnitude than regression coefficients later in the sequence and subtract the number of discordant pairs. The decision rule then predicts the direction in which the sequence of regression coefficients is more increasing according to this criterion.

We apply this orientation rule to simulated data (sample size 1000) for varying chain lengths and edge distributions, and when applied to raw observational data, standardized observational data, and data when the parameters were scale-harmonized as per Mooij et al. [2020]. The table below establishes, that for iid distributed parameters of the underlying data generating process, the orientation of a causal chain can be identified with probability strictly greater than 50%.

Table 4: Empirical Chain Orientation Results

d	edge range	accuracy by variance-sorting			accuracy by coefficient-sorting		
		raw	standardized	harmonized	raw	standardized	harmonized
3	$\pm(0.5, 2.0)$	97.50%	50.05%	84.70%	62.58%	73.03%	57.30%
	$\pm(0.5, 0.9)$	80.38%	50.05%	69.62%	57.15%	62.38%	55.65%
	$\pm(0.1, 0.9)$	65.65%	50.30%	60.08%	54.17%	55.88%	53.45%
5	$\pm(0.5, 2.0)$	98.67%	50.15%	82.17%	78.60%	86.58%	64.20%
	$\pm(0.5, 0.9)$	77.65%	49.27%	66.30%	61.83%	68.65%	57.50%
	$\pm(0.1, 0.9)$	63.08%	50.38%	57.65%	58.17%	57.33%	56.35%
10	$\pm(0.5, 2.0)$	99.38%	50.02%	79.30%	93.72%	96.97%	69.08%
	$\pm(0.5, 0.9)$	73.75%	50.25%	62.00%	64.97%	70.70%	58.50%
	$\pm(0.1, 0.9)$	62.55%	51.23%	58.25%	55.85%	56.05%	54.40%

A 99% confidence interval for the orientation accuracy under random guessing is (47.975%, 52.025%) (1000 repetitions for each of the four noise types). Thus, variance-sorting on the standardized data is the only setting in which no above-chance orientation accuracy is achieved. This is expected, as variance sorting amounts to a random sorting once nodes are standardized.

G Empirical Evaluation of Varsortability

We empirically estimate expected varsortability for our experimental set-up and a non-linear version of our experimental set-up by calculating the fraction of directed paths that are correctly sorted by marginal variance in the randomly sampled ANMs.

G.1 Varsortability in Linear Additive Noise Models

Consistent with our theoretical results, varsortability is close to 1 across all graph and noise types in our experimental set-up, cf. Table 5. Varsortability is higher in denser than in sparser graphs.

Table 5: Empirical varsortability in our experimental linear ANM set-up. Average varsortability is high in all settings. Our parameter choices are common in the literature. We sample 1000 observations of ten 50-node graphs for each combination of graph and noise type.

graph	noise	varsortability		
		min	mean	max
ER-1	Gauss-EV	0.94	0.97	0.99
	exponential	0.94	0.97	0.99
	gumbel	0.94	0.97	1.00
ER-2	Gauss-EV	0.97	0.99	1.00
	exponential	0.97	0.99	1.00
	gumbel	0.98	0.99	0.99
ER-4	Gauss-EV	0.98	0.99	0.99
	exponential	0.98	0.99	0.99
	gumbel	0.98	0.99	0.99
SF-4	Gauss-EV	0.98	1.00	1.00
	exponential	0.98	1.00	1.00
	gumbel	0.98	1.00	1.00

G.2 Varsortability in Non-Linear Additive Noise Models

Table 6 shows varsortabilities for a non-linear version of our experimental set-up as used by [Zheng et al. \[2020\]](#). While the fluctuations in Table 6 are greater than in Table 5, all settings exhibit high varsortability on average. Our findings indicate that varsortability is a concern for linear and non-linear ANMs.

Table 6: Empirical varsortability in non-linear ANM. Average varsortability is high in all settings. Our parameter choices are common in the literature. We sample 1000 observations of ten 20-node graphs for each combination of graph and ANM-type.

graph	ANM-type	varsortability		
		min	mean	max
ER-1	Additive GP	0.81	0.91	1.00
	GP	0.72	0.86	0.96
	MLP	0.55	0.79	0.96
	Multi Index Model	0.62	0.82	1.00
ER-2	Additive GP	0.79	0.91	0.98
	GP	0.82	0.89	0.97
	MLP	0.46	0.71	0.87
	Multi Index Model	0.65	0.79	0.89
ER-4	Additive GP	0.90	0.95	0.98
	GP	0.74	0.88	0.93
	MLP	0.59	0.72	0.85
	Multi Index Model	0.57	0.73	0.85
SF-4	Additive GP	0.95	0.97	0.99
	GP	0.88	0.94	0.97
	MLP	0.75	0.83	0.93
	Multi Index Model	0.77	0.84	0.97

G.3 Causal Order and Marginal Variance

We observe strong empirical evidence in Figure 2 that marginal variance tends to increase quickly along the causal order, even if the settings are not guaranteed to yield high expected varsortability between a pair of root cause and effect (for example, if all edges are chosen in a small-magnitude range). This indicates that high levels of varsortability can scarcely be avoided on larger graphs.

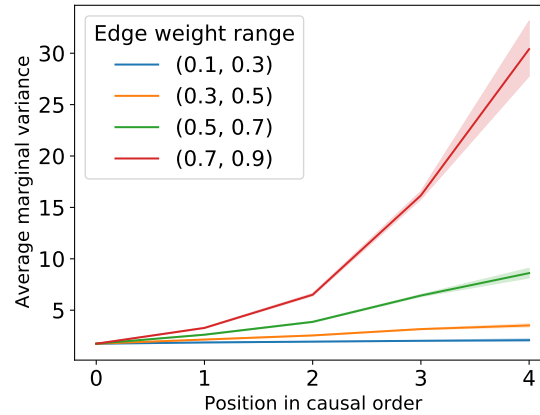


Figure 2: Average marginal variance along the causal order for 1000 observations of 1000 simulated 30-node ER-2 graphs with Gaussian noise standard deviations sampled uniformly in $(0.5, 2)$ for each edge weight range. Edge weights are drawn independently and uniformly from the union of negative and positive of the indicated edge range, that is, for example, the edge weights for the red curve are drawn from $\text{Unif}((-0.9, -0.7) \cup (0.7, 0.9))$.

G.4 Varsortability Algorithm

The implementation is also available at <https://github.com/Scriddie/Varsortability>.

```
import numpy as np
```

```
def varsortability(X, W, tol=1e-9):
    """ Takes  $n \times d$  data and a  $d \times d$  adjacency matrix,
    where the  $i, j$ -th entry corresponds to the edge weight for  $i \rightarrow j$ ,
    and returns a value indicating how well the variance order
    reflects the causal order. """
    E = W != 0
    Ek = E.copy()
    var = np.var(X, axis=0, keepdims=True)

    n_paths = 0
    n_correctly_ordered_paths = 0

    for _ in range(E.shape[0] - 1):
        n_paths += Ek.sum()
        n_correctly_ordered_paths += (Ek * var / var.T > 1 + tol).sum()
        n_correctly_ordered_paths += 1/2*(
            (Ek * var / var.T <= 1 + tol) *
            (Ek * var / var.T > 1 - tol)).sum()
        Ek = Ek.dot(E)

    return n_correctly_ordered_paths / n_paths

if __name__ == "__main__":
    W = np.array([[0, 1, 0], [0, 0, 2], [0, 0, 0]])
    X = np.random.randn(1000, 3).dot(np.linalg.inv(np.eye(3) - W))
    print("Varsortability:", varsortability(X, W))

    X_std = (X - np.mean(X, axis=0))/np.std(X, axis=0)
    print("Varsortability_σstandardized:", varsortability(X_std, W))
```

H *sortnregress*: A Diagnostic Tool to Reveal Varsortability

In Section 3.5 we introduce *sortnregress* as a simple baseline method. In the following subsections, we provide Python code that implements *sortnregress* thereby establishing its ease and illustrate how its DAG recovery performance reflects varying degrees of varsortability.

H.1 Implementation of Sortnregress

The implementation is also available at <https://github.com/Scriddie/Varsortability>.

```
import numpy as np
from sklearn.linear_model import LinearRegression, LassoLarsIC

def sortnregress(X):
    """ Take n x d data, order nodes by marginal variance and
    regresses each node onto those with lower variance, using
    edge coefficients as structure estimates. """
    LR = LinearRegression()
    LL = LassoLarsIC(criterion='bic')

    d = X.shape[1]
    W = np.zeros((d, d))
    increasing = np.argsort(np.var(X, axis=0))

    for k in range(1, d):
        covariates = increasing[:k]
        target = increasing[k]

        LR.fit(X[:, covariates], X[:, target].ravel())
        weight = np.abs(LR.coef_)
        LL.fit(X[:, covariates] * weight, X[:, target].ravel())
        W[covariates, target] = LL.coef_ * weight

    return W
```

H.2 Varsortability and Score Attainable by Variance Ordering

In Figure 3 we observe that *sortnregress* improves linearly with varsortability. For a varsortability of 0.93 as in our experimental settings (cf. Section 3.3), it recovers the structure near-perfectly. *randomregress* uses a random ordering but is otherwise identical to *sortnregress*. The different ranges of varsortability can be classified as follows (n=30):

- < 0.33: *sortnregress* performs significantly worse than *randomregress* (p<1e-4)
- 0.33–0.66: no significant difference between *sortnregress* and *randomregress* (p=0.40)
- > 0.66: *sortnregress* performs significantly better than *randomregress* (p<1e-4)

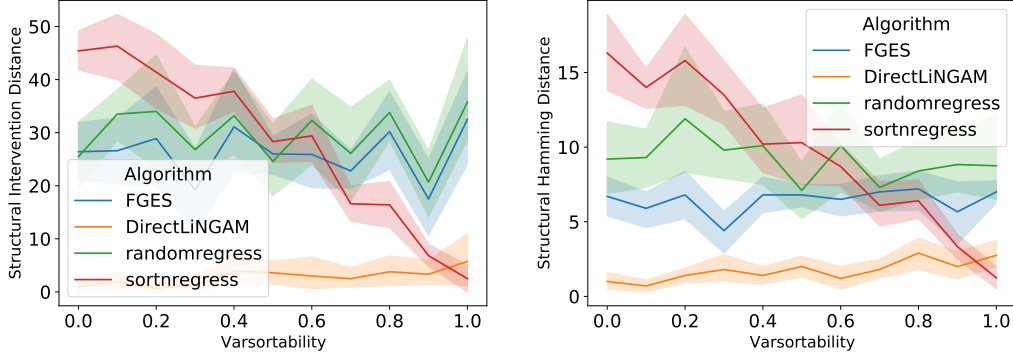


Figure 3: Relationship between varsortability and score attainable through ordering by variance. Results shown for 10 simulated 10-node ER-1 graphs in each of 10 equally spaced varsortability bins. Note that for standard simulation settings most models have high varsortability. We use edge weights in $(-0.5, -0.1) \cup (0.1, 0.5)$, Gumbel noise with standard deviations in $(0.5, 2)$, and still need to discard many models with high varsortability to obtain 10 instances per varsortability bin.

I Evaluation on Real-World Data

We analyze a dataset on protein signaling networks obtained by [Sachs et al. \[2005\]](#). We evaluate our algorithms on ten bootstrap samples of the observational part of the dataset consisting of 853 observations, 11 nodes, and 17 edges. Our results show that there is no dominating algorithm. On average, most algorithms achieve performances similar to those of *randomregress* or the empty graph. Note that the results in terms of SHD are susceptible to thresholding choices and the empty graph baseline outperforms a majority of the algorithms. Our results are in line with previous reports [[Lachapelle et al., 2019](#), [Ng et al., 2020](#)]. We observe scale-sensitivity of the continuous learning algorithms and *sortnregress*. However, in contrast to our simulation study in Section 4, the effect is small and inconsistent. The results do not show the patterns observed under high varsortability, which is consistent with the measured mean varsortability of 0.57 with a standard deviation of 0.01 across our bootstrapped samples.

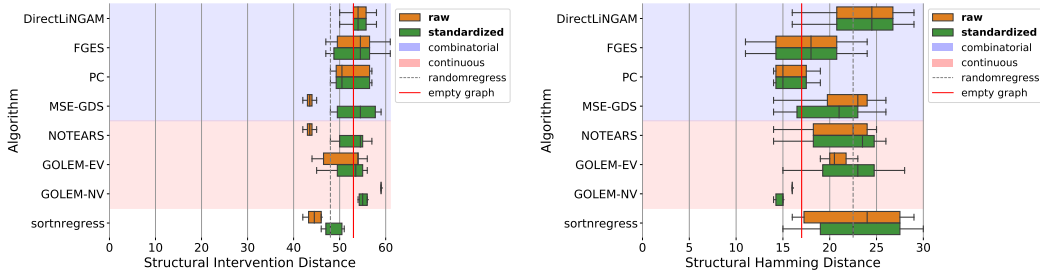


Figure 4: SID (left) and SHD (right) performance of combinatorial and continuous methods on real-world data.

J Model Selection in Continuous Optimization

We illustrate the optimization landscape for the Gaussian MLE under Gaussian noise. This corresponds to the loss of *GOLEM-NV* as stated in Appendix D with a sparsity penalty of zero. We compare vanilla MLE to MLE with Lasso regularization for raw and standardized data. In Figure 5 we show the loss landscape in terms of SID and SHD difference to the true structure and highlight global optima. In the case of tied scores between the true structure and an alternative structure we select the true structure. For MLE with Lasso regularization using a penalty of 0.1, the optimal loss is achieved by the true structure more frequently under standardization (red dots accumulate in the

bottom left corner). Our result indicates that the Lasso sparsity penalty is influenced by the data scale and is better calibrated on standardized data. It is not unexpected that penalization is scale dependent, a problem that is, for example, discussed in applications of Ridge regression.

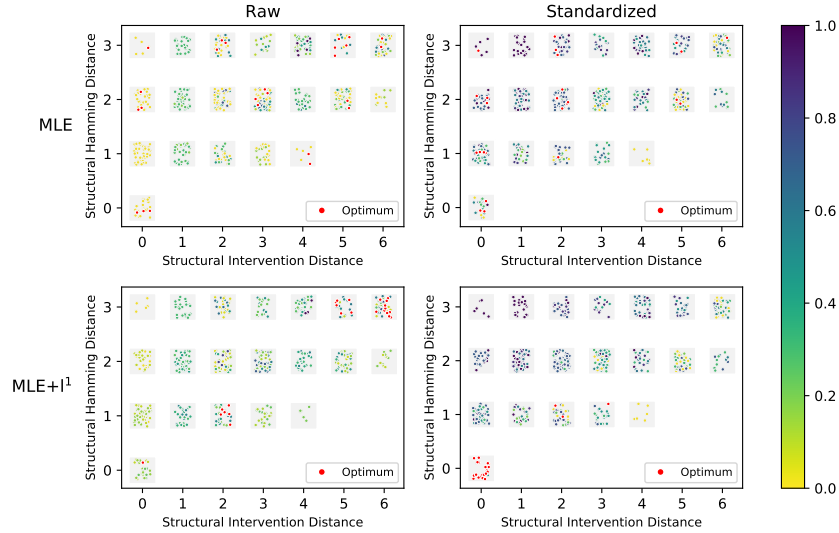


Figure 5: Standardized loss landscape for all 25 candidate graphs relative to each of the 25 possible 3-node ground-truth structures (a total of 25×25 candidate-true graph pairs). The loss is scaled to $[0, 1]$, see colorbar.

K Detailed Results

We provide a comprehensive overview over our empirical DAG/MEC recovery results for different evaluation metrics, graph types, and graph sizes.

K.1 MEC Recovery

An analysis of MEC recovery allows us to distinguish whether any drops in performance are within the expectations of identifiability. We evaluate the discovery of the MEC of the ground-truth DAG in a Gaussian setting with non-equal noise variances where only the ground-truth MEC but not the ground-truth DAG are identifiable. Since evaluating the SID between Markov equivalence classes is computationally expensive and prohibitively so for large graphs, we restrict ourselves to the setting here. When comparing MEC, we choose the upper limit of SID differences in Figure 1 in the main text. In Figure 6 we show that the relative performances are similar for the lower SID limit.

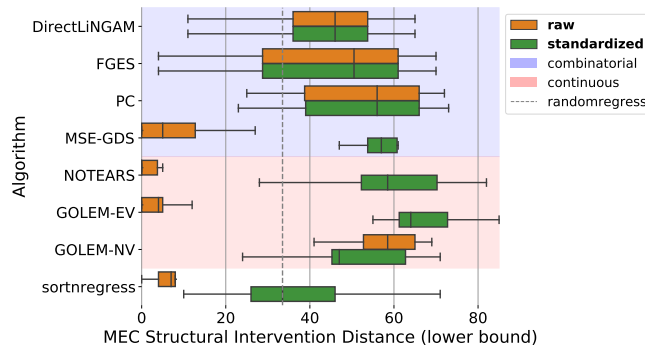


Figure 6: Lower bound of SID in MEC recovery for 10 node ER-2 graphs with non-equal Gaussian noise.

We conclude that the drop in performance extends from the recovery of the DAG to the recovery of the MEC and therefore goes beyond the difficulty of identifying the correct DAG within a MEC.

K.2 Results Across Thresholding Regimes

To ensure the effects we observe constitute a general phenomenon, we evaluate algorithm performance for different thresholding regimes. This is especially critical on standardized data. By re-scaling the data, standardization may impact the correct edge weights between nodes, potentially pushing them outside the thresholding range. Following Zheng et al. [2018], Ng et al. [2020], we perform thresholding for the continuous structure learning algorithms and prune edges with an edge weight in the recovered adjacency matrix of less than 0.3. If the returned graph is not acyclic, we iteratively remove the edge with the smallest magnitude weight until all cycles are broken. We find that the qualitative performance differences between raw and standardized data are robust to a wide range of threshold choices.

Figure 7a and Figure 7b show SID performance for different thresholds. Even though the thresholds are orders of magnitude apart, a comparison reveals that the relative performances are nearly identical.

We observe that SHD performance is also robust across different thresholding regimes. Figure 7c shows performance using *favorable* thresholding. In this regime, the threshold leading to the most favorable SHD performance is applied to each instance individually. Figure 7d shows performance for a fixed threshold of 0.3. A comparison reveals nearly identical relative performances in both cases.

Overall, we observe that the effect of varsortability is present even for the most favorable threshold in case of SHD, and for a wide range of thresholds in case of SID, where computation of a favorable threshold is computationally infeasible.

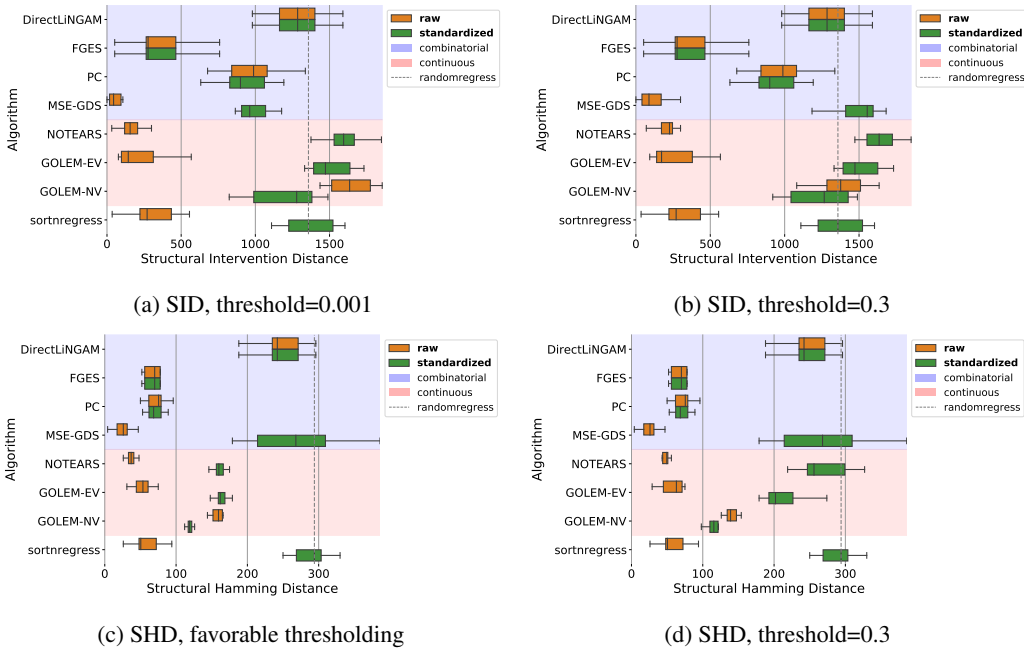
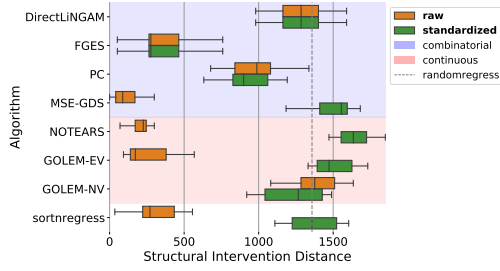


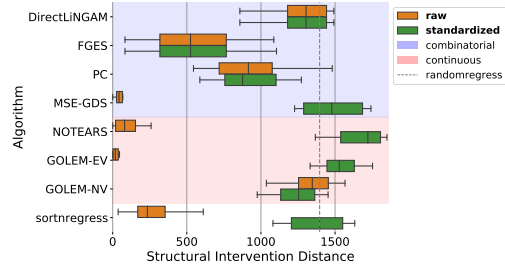
Figure 7: Results for different thresholding regimes. Gaussian-NV noise, ER-2 graph, 50 nodes.

K.3 Results Across Noise Distributions and Graph Types

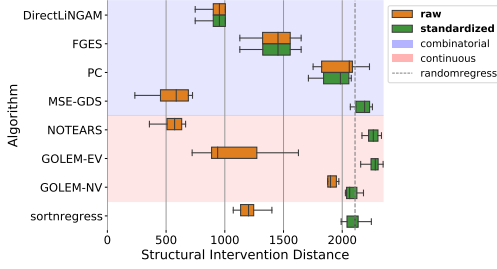
Figure 8 and Figure 9 show algorithm comparisons in terms of SID and SHD, respectively. The differences in performance on raw versus standardized data are qualitatively similar regardless of the noise distribution. We showcase results for different graph types in the non-Gaussian setting. *DirectLINGAM* performs well only in the non-Gaussian cases, as is expected based on its underlying identifiability assumptions.



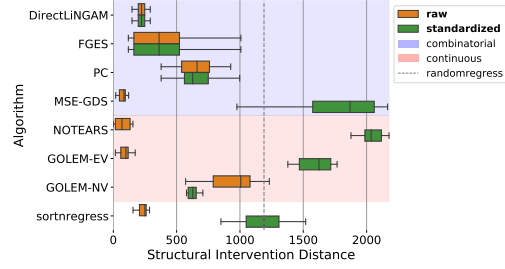
(a) SID, Gaussian-NV noise, ER-2 graph, 50 nodes



(b) SID, Gaussian-EV noise, ER-2 graph, 50 nodes

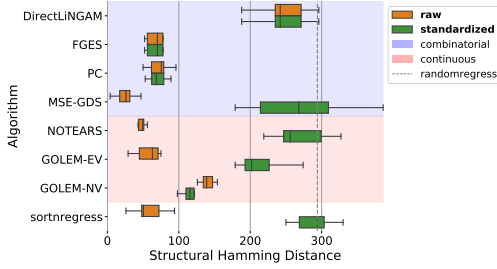


(c) SID, Exponential noise, ER-4 graph, 50 nodes

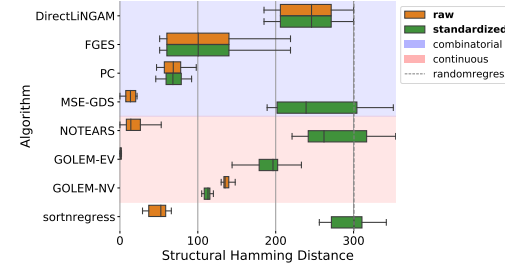


(d) SID, Gumbel noise, SF-4 graph, 50 nodes

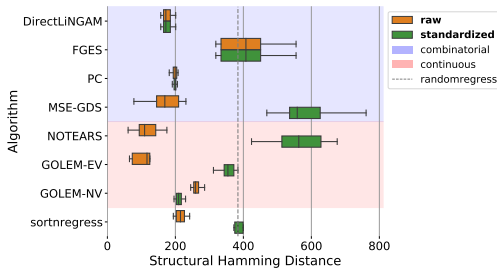
Figure 8: SID results across noise types and for different graph types with 50 nodes



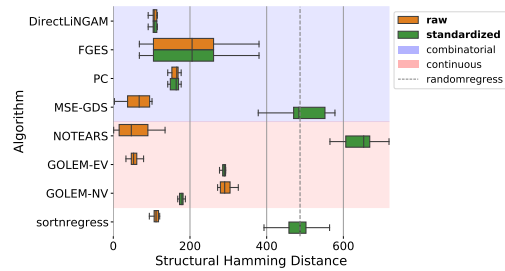
(a) SHD, Gaussian-NV noise, ER-2, 50 nodes



(b) SHD, Gaussian-EV noise, ER-2, 50 nodes



(c) SHD, Exponential Noise, ER-4, 50 nodes



(d) SHD, Gumbel Noise, SF-4, 50 nodes

Figure 9: SHD results across noise types and for different graph types with 50 nodes

K.4 Results Across Noise Distributions, Graph Types, and Graph Sizes

The following experimental results largely follow earlier settings and results by [Zheng et al. \[2018\]](#), [Ng et al. \[2020\]](#).

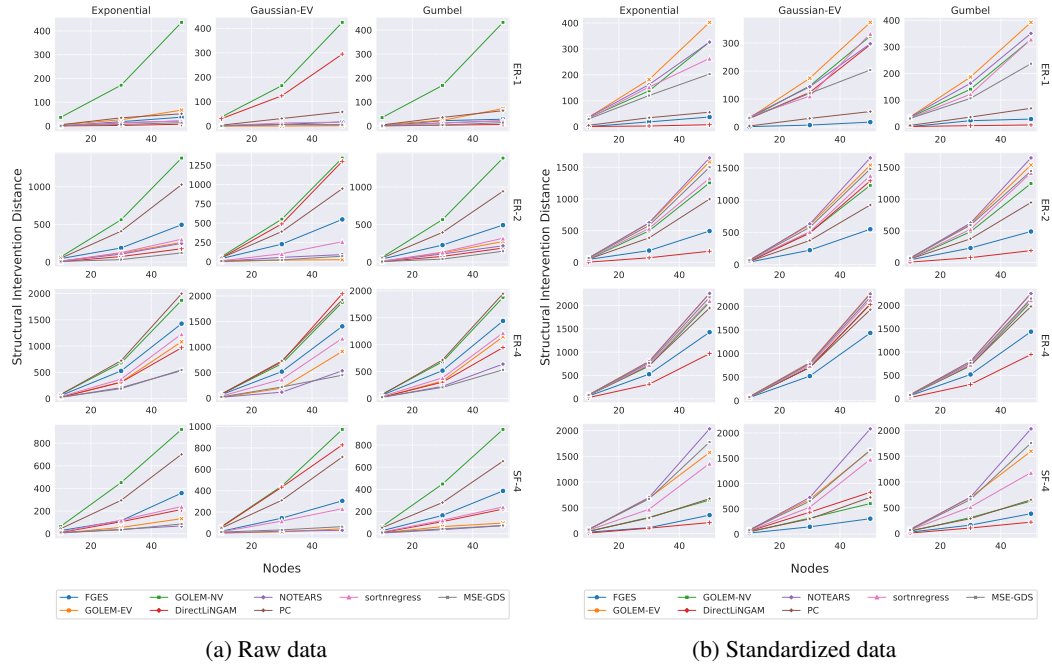


Figure 10: SID results across noise types, graph types, and graph sizes.

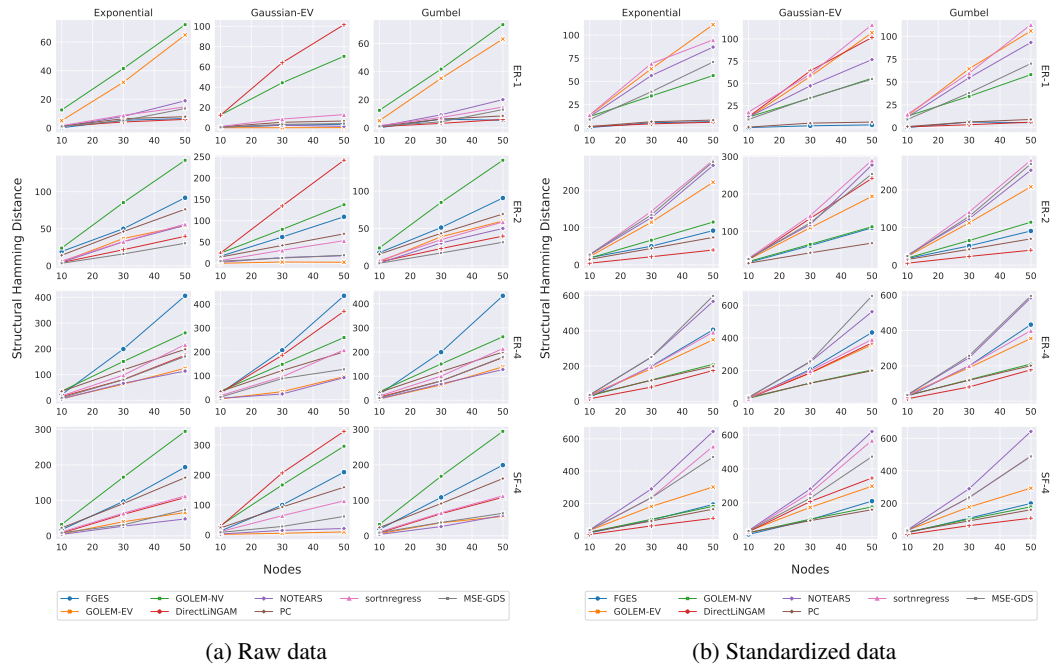


Figure 11: SHD results across noise types, graph types, and graph sizes.