
Supplementary Material

Anonymous Author(s)
Affiliation
Address
email

1 Example DAGs

2 In order to evaluate the different causal discovery methods we generated random Directed Acyclic
3 Graphs (DAGs) using eight different generation methods of the randDAG function of the R pcalg
4 library¹. These generation methods were *regular* - a graph where every node has exactly d incident
5 edges, *er* - an Erdos-Renyi graph where every edge is present independently, *watts* - an interpolation
6 between regular graph and Erdos-Renyi graph, *power* - a graph with power-law degree distribution,
7 *bipartite* - a bipartite graph, *barabasi* - a graph with power-law degree distribution and preferential
8 attachment, *geometric* - a geometric random graph, and *interEr* - a graph with two islands of Erdos-
9 Renyi graphs connected by a small number of edges. An example of a graph produced by these
10 different methods can be seen in Figure 1.

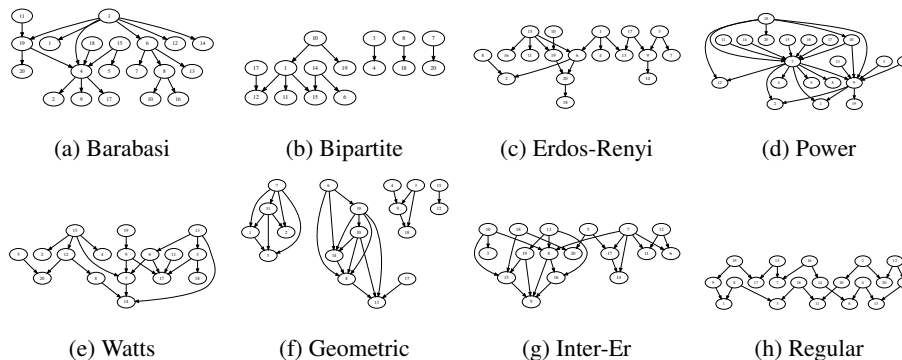


Figure 1: Examples of causal graphs with $N = 20$ nodes and $E = 2$ expected edges per node, generated using each method of the PCALG random DAG function. We generated 1140 graphs total with $N = 20, N = 40, N = 60$ nodes and $E = 1, E = 2, E = 3, E = 4,$ and $E = 5$ expected edges per node.

11 2 Stability Algorithm Selection Heuristic

12 We focus on detailed evaluation of the agreement algorithm selection heuristic in the main paper, but
13 we also explored an alternative heuristic of algorithm stability, or robustness to sampling of the data,
14 as another potential indicator of algorithm accuracy. We measure the stability of a given algorithm by
15 applying the algorithm to multiple samples of the data. For each edge or node that appears in any
16 graph generated from these samples, we calculated the percentage of the graphs in which the given
17 edge or node appears. An edge which appears in all graphs is highly stable, while one that appears
18 only in one is not. This definition of stability focuses on the consistency of edge presence rather

¹<https://www.rdocumentation.org/packages/pcalg/versions/2.6-8/topics/randDAG>

19 than the consistency of edge absence. To calculate the overall stability of the graph, we average the
 20 stability of each individual edge or node.

21 In Figure 2, we show that algorithm stability is positively correlated with the accuracy of the algorithm.
 22 In fact, the correlation is observed to be positive for all graph generation methods and parameters.
 23 However, the level of correlation is significantly less on average than is observed for the agreement
 24 scores. Therefore, we expect that the agreement score will provide a better indicator for accuracy and
 25 incorporate this score in both the algorithm selection and ensemble methods that we develop.

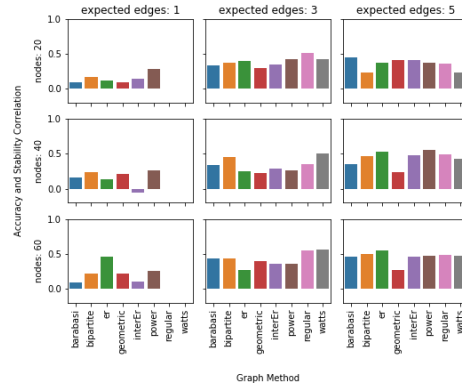


Figure 2: The Spearman rank correlation of algorithm accuracy and algorithm stability for PCALG data generated from graphs with different structural properties.

26 3 Parameter Optimization

27 The best performing ensemble method used the subselected four algorithms, no upweighting of
 28 agreeable algorithms, and an edge weight threshold of 0.65. In Figure 3, we show the a comparison
 29 of the ensemble performance across the algorithm selection and upweighting ratio parameters for
 30 datasets generated from graphs with different structural properties. We find that for the majority of
 31 graph types, the four-algorithm ensemble outperforms the full ensemble for the full range of ratio
 32 settings. This indicates that the performance of the ensemble approach is sensitive to the selection of
 33 specific algorithms to include, with lower performing algorithms reducing the overall performance
 34 of the ensemble. We also find that the two different algorithm selections have different patterns with
 35 respect to the upweighting ratio. We find that increasing the weight of the most agreeable algorithm
 36 increases the performance of the full ensemble. In contrast, for the four-algorithm ensemble, we find
 37 that upweighting the most agreeable algorithms reduces the performance. This difference is likely
 38 due to two factors. Firstly, with fewer algorithms included the overall agreement score may become
 39 less meaningful because there are fewer inter-algorithm comparison points. Secondly, because the
 40 agreement scores are correlated with accuracy, the agreement score upweighting in the full ensemble
 41 is effectively downweighting the less accurate algorithms. These algorithms are excluded entirely
 42 from the selective ensemble, so the downweighting becomes less necessary.

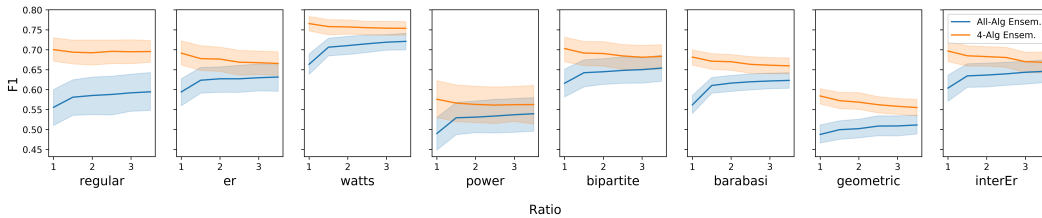


Figure 3: Causal discovery performance for the all-algorithm and the four-algorithm ensembles as a function of how much the most agreeable algorithm is upweighted (the ratio of the most agreeable algorithm weight to the least agreeable weight) evaluated on PCALG data.

