

An Algorithm for Fast Recovery of Sparse Causal Graphs

by

Peter Spirtes, Clark Glymour

August 1990

Report CMU-PHIL-15



Philosophy
Methodology
Logic

Pittsburgh, Pennsylvania 15213-3890

An Algorithm for Fast Recovery of Sparse Causal Graphs

Peter Spirtes and Clark Glymour¹
Carnegie Mellon University

Abstract

Previous asymptotically correct algorithms for recovering causal structure from sample probabilities have been limited even in sparse graphs to a few variables. We describe an asymptotically correct algorithm whose complexity for fixed graph connectivity increases polynomially in the number of vertices, and may in practice recover sparse graphs with several hundred variables. From sample data with $n = 2,000$, an implementation of the algorithm on a Decstation 3100 recovers the edges in a linear version of the ALARM network with 37 vertices and 46 edges. Fewer than 8% of the undirected edges are incorrectly identified in the output. Without prior ordering information the program also determines the direction of edges for the ALARM graph with an error rate of 14%. Processing time is less than 10 seconds.

Consider pairs $\langle g, P \rangle$ for which g is a directed acyclic graph and P is a probability distribution on the vertices of g such that (i) for every vertex v and every set S_v of vertices that are not descendants of v , v and S_v are independent conditional on the parents of v ; and (ii) every independence relation in P is a consequence of the independence relations in (i). Pairs satisfying these conditions can be viewed as causal structures in which the causal dependencies generate statistical dependencies. When the set of measured variables for which probabilities are provided in the data is such that every common cause of a measured variable is measured, we say the structure is causally sufficient.

¹We thank Gregory Cooper for a conversation that stimulated this work.

1 August, 1990

Recovery problems have to do with determining g , or features of g , from the distribution P or from samples obtained from P . In Spirtes [1990a, 1990b] we proposed the following algorithm for the recovery problem with causally sufficient structures, using as input independence and conditional independence facts from P^2 :

SGS Algorithm

(1) For each vertex pair a, b , and each subset S of vertices not containing a or b , determine if $I(a, S, b)$. If such an S is found go to the next pair of vertices; if no such vertex is found, place an undirected edge between a and b , and go to the next pair of vertices until all pairs have been considered.

(2) For each triple a, b, c of vertices such that a and b are adjacent, b and c are adjacent and a and c are adjacent, direct the edges $a - b$ and $b - c$ into b if and only if for every set S of vertices containing b but not a or c , $\neg I(a, S, b)$

(3) Output all orientations of the graph consistent with (2).

Verma and Pearl [Verma (1990)] subsequently proved the correctness of the algorithm and offered a variant that outputs a pattern rather than a collection of graphs. The pattern has an undirected edge between two vertices if the SGS output contains graphs that orient the edge in both directions; the pattern contains a directed edge if every graph output by the SGS algorithm has the edge so oriented; and the pattern may have a bidirected edge, e.g., $a \leftrightarrow b$ provided step two of the algorithm determines that the $a - b$ edge collides with another edge at a and also collides with another edge at b . When all common causes are measured and the data consist of the actual independence and conditional independence relations, the pattern is

²We denote by " $I(a, S, b)$ " the claim that variables a and b are independent conditional on the set of variables in S , and by " $\neg I(a, S, b)$ " the denial of that claim.

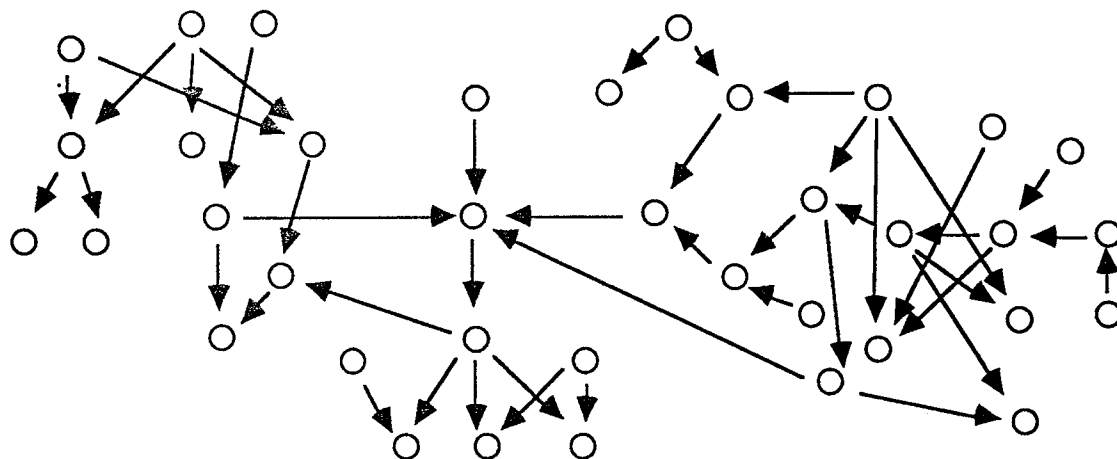
simply a representation of the class output by the SGS algorithm, but when there are unmeasured common causes or independence facts due to sampling variation rather than to P , the pattern is more general.

Two graphs, g, g' are statistically indistinguishable provided that for every probability distribution P , $\langle g, P \rangle$ satisfies the conditions (i) and (ii) of the first paragraph if and only if $\langle g', P \rangle$ does. From the independence facts of a distribution P such that $\langle g, P \rangle$ satisfies (i), and (ii), the SGS algorithm returns all and only the graphs statistically indistinguishable from g .

In the worst case, the SGS algorithm requires a number of conditional independence facts that increases exponentially with the number of vertices, as must any algorithm based on conditional independence relations. But because the number of conditional independence facts that must be generated and checked in stage (1) of the algorithm is unaffected by the connectivity of the true graph, even for sparse graphs the algorithm rapidly becomes computationally infeasible as the number of vertices increases. Besides problems of computational feasibility, the algorithm has problems of reliability when applied to sample data. The determination of higher order conditional independence relations from sample distributions is generally less reliable than is the determination of lower order independence relations. With, say, 37 binary variables, to determine the conditional independence of two variables on the set of all remaining variables requires considering the relations among the frequencies of 2^{35} distinct states, only a tiny fraction of which will be instantiated even in very large samples.

To illustrate the difficulty consider an example due to Herskovits and Cooper (1990). Their Kutato' Algorithm is a heuristic entropy minimization procedure for recovering a directed graph given sample data and a total ordering of the vertices such that $v_1 > v_2$ implies that there is no directed edge from v_2 to v_1 . The asymptotic

reliability of the procedure is unknown. Nonetheless from large sample data the algorithm recovers most of the connections on a sparse graph--the ALARM network (Beinlich, 1989)--with 37 variables and 46 edges. In their example, the direction of the edges is not recovered from the data but is determined by the prior ordering given to the computer.³



ALARM Network

Using 10,000 cases an implementation on a MacIntosh II required about 22 and one half hours, about a quarter of which was required to read the data-base. The output omitted two correct edges and included two false edges. By comparison, the SGS algorithm has been implemented in the TETRAD II program using partial correlation tests for conditional independence. Run on a DEC workstation with 20 megabyte RAM the procedure stops at about 17 variables because

³Herskovits and Cooper say that a variant of the Kutato' algorithm can determine the orientation of edges without a prior ordering of the variables, but they do not describe the properties of the application or give an example. They are also investigating Bayesian alternatives that they expect to be faster than the Kutato' procedure.

PC Algorithm:

Let A_{ab} denote the set of vertices adjacent to a or to b . Let P_{ab} denote the set of vertices that are parents of a or of b . Let U_{ab} denote the set of vertices on undirected paths between a and b .

- A.) Form the complete graph C_{-1} on the vertex set V .
- B.) For each pair of variables a, b adjacent in C_n :
 - (α) If $A_{ab} \cap U_{ab}$ does not have cardinality greater than n , go to the next pair of vertices adjacent in C_n .
 - (β) If A_{ab} has does have cardinality greater than n , determine if a, b are independent conditional on any subsets of $A_{ab} \cap U_{ab}$ of cardinality $n + 1$. If so, delete $a-b$ from C_n .

Let C_{n+1} be the graph that results from this procedure applied to each pair of variables. Continue until a value $f + 1$ of n is reached such that β is not satisfied for any pair.

- C.) For each triple of vertices a, b, c such that the pair a, b and the pair b, c are each adjacent in C_f but the pair a, c are not adjacent in C_f , orient $a - b - c$ as $a \rightarrow b \leftarrow c$ if and only if a and c are dependent on every subset of $A_{ac} \cap U_{ac}$ containing b .
Output all graphs consistent with these orientations.

Note that $A_{ab} \cap U_{ac}$ is not in general the set of *parents* of a or b on undirected paths between a, b , since descendants of a, b may also occur.

An obvious modification of the algorithm will generate patterns rather than collections of graphs.

The complexity of the algorithm for a graph G is bounded by $\max(|A_{ab}|)$ over all pairs of vertices a, b , which is never more than the sum of the two largest degrees in G . Generally stage B of the

of space requirements to store the conditional independence facts. Space could be traded for time, but the ALARM case is out of sight.

Verma and Pearl have suggested an improvement on the SGS algorithm. For each pair of variables a, b introduce an undirected edge between them if they are dependent conditional on the set of all other variables. Call the resulting network N . (In the true graph, G , the parents of any variable form a maximal complete subgraph--a clique--in the network N .) Again for each pair of variables a, b adjacent in N , determine if a, b are dependent conditional on all subsets of variables in the cliques in N containing a or b . If so a is adjacent to b in G . The complexity is thus bounded by the size of the largest clique in N .

The practical value of the improvement is limited by the fact that conditional independence relations of the order of the number of vertices of the graph (minus two) must still be estimated, with consequent costs in computational efficiency and reliability. With discrete data the great majority of the corresponding states will not be instantiated in the data, and with data from linear structures the formula for higher order correlations is recursive: to compute the partial correlations of n th order, three partial correlations of order $n-1$ must be determined, and so on.

We should like an algorithm that has the same input/output relations as the SGS procedure but for sparse graphs does not require the determination of higher order independence relations, and in any case requires as few conditional independence relations as possible. The following procedure starts by forming the complete undirected graph, then "thins" that graph by removing edges with zero order conditional independence relations, thins again with first order conditional independence relations, and so on. The set of variables conditioned on need only be a subset of the set of variables adjacent to one or the other of the variables conditioned, and can even be confined to adjacent variables on certain undirected paths.

PC Algorithm:

Let A_{ab} denote the set of vertices adjacent to a or to b . Let P_{ab} denote the set of vertices that are parents of a or of b . Let U_{ab} denote the set of vertices on undirected paths between a and b .

- A.) Form the complete graph C_1 on the vertex set V .
- B.) For each pair of variables a, b adjacent in C_n :
 - (α) If $A_{ab} \cap U_{ab}$ does not have cardinality greater than n , go to the next pair of vertices adjacent in C_n .
 - (β) If A_{ab} has does have cardinality greater than n , determine if a, b are independent conditional on any subsets of $A_{ab} \cap U_{ab}$ of cardinality $n + 1$. If so, delete $a-b$ from C_n .

Let C_{n+1} be the graph that results from this procedure applied to each pair of variables. Continue until a value $f + 1$ of n is reached such that β is not satisfied for any pair.

- C.) For each triple of vertices a, b, c such that the pair a, b and the pair b, c are each adjacent in C_f but the pair a, c are not adjacent in C_f , orient $a - b - c$ as $a \rightarrow b \leftarrow c$ if and only if a and c are dependent on every subset of $A_{ac} \cap U_{ac}$ containing b . Output all graphs consistent with these orientations.

Note that $A_{ab} \cap U_{ac}$ is not in general the set of *parents* of a or b on undirected paths between a, b , since descendants of a, b may also occur.

An obvious modification of the algorithm will generate patterns rather than collections of graphs.

The complexity of the algorithm for a graph G is bounded by $\max(|A_{ab}|)$ over all pairs of vertices a, b , which is never more than the sum of the two largest degrees in G . Generally stage B of the

algorithm continues testing for some steps after the correct undirected graph has been identified. The number of steps required before the true graph is found (but not necessarily until the algorithm halts) depends on the maximal number of treks between a pair of variables, say a , b , such that no two treks⁴ share vertices adjacent to a or b . If these maximal numbers are held constant as the number of vertices increases, so that k , the maximal order of the conditional independence relations that need be tested, does not change, then the worst case computational demands of the algorithm increase as

$$\frac{n!}{2!(n-2)!} \sum_{j=2}^k \frac{(n-2)!}{j!(n-2-j)!}$$

which is bounded by n^k . It should be possible to recover sparse graphs with as many as several hundred variables. Of course the computational requirements increase exponentially with k .

In many cases it may be more efficient to perform conditional independence tests on all subsets of the set of adjacent variables rather than to compute the sets of adjacent variables that lie on undirected paths. We have not yet theoretically determined the trade-off.

The structure of the algorithm and the fact that it continues to test even after having found the correct graph suggest a natural heuristic for very large variable sets whose causal connections are expected to be sparse, namely to set a fixed bound on the order of conditional independence relations that will be considered.

Proposition: The PC and SGS algorithms give the same output.

⁴A trek is a pair of directed paths from some vertex z to a , b respectively, intersecting only at z , or a directed path from a to b or a directed path from b to a .

Proof:

We note a lemma:

Lemma: In any pair $\langle g, P \rangle$ meeting conditions i, and ii. if vertices a, b are not adjacent then they are independent conditional on $P_{ab} \cap U_{ab}$.

The proof is a trivial modification of the argument Verma and Pearl give for their Lemma 1.

Now we show that steps A and B of the PC algorithm produce the correct undirected graph. The undirected graph of G is a subgraph of all graphs $C_1 \dots C_f$. Two variables, a, b are adjacent in G if and only if they are dependent on every set of vertices containing neither of them. Trivially, if variables a and b are dependent conditional on each set of variables not containing them, then they are dependent conditional on every subset of $A_{ab} \cap U_{ab}$. We must show the converse: If a, b are dependent conditional on every subset of $A_{ab} \cap U_{ab}$ then a, b are adjacent.

Suppose a, b are not adjacent. Then by the Lemma a, b are independent conditional on the set $P_{ab} \cap U_{ba}$ which is contained in $A_{ab} \cap U_{ab}$.

It remains only to show that step C of the algorithm orients the graph correctly. Assume a, c are not adjacent but a is adjacent to b and b is adjacent to c . It is known that the $a - b$ and $b - c$ edges collide at b if and only if there is no set S containing b and not a or c such that a, c are independent conditional on S . Since a, c are not adjacent, they are independent conditional on the set $P_{ac} \cap U_{ac}$. If the edges do not collide at b , then b is a parent of a or of c , so b is in $P_{ac} \cap U_{ac}$. If the edges do collide at b then a, c are dependent on

every set containing b and not a , c , and hence dependent on every subset of $A_{ac} \cap U_{ac}$ that contains b .

We have applied the PC algorithm to a linear version of the ALARM network. Using the same directed graph, linear coefficients with values between .5 and 1.0 were randomly assigned to each directed edge in the graph. Using a joint normal distribution on the variables of zero indegree, three sets of simulated data were generated, each with a sample size of 2,000. The covariance matrix and sample size were given to a version of the TETRAD II program with an implementation of the PC algorithm that does not check to determine whether variables adjacent to vertices v_1 , v_2 lie on an undirected path between v_1 and v_2 , and that outputs a pattern. No information about the orientation of the variables was given to the program. Run on a Decstation 3100, for each data set the program required less than ten seconds to return a pattern. In each trial the output pattern omitted three edges in the ALARM network. Of the remaining 43 edges, the orientation of three of them is not determinable in principle from the probabilities, and in the first two trials the program so reported while in the third it oriented one of the three. Of the remaining forty edges, the trials misoriented 5, 6, and 7 edges respectively, always by judging that an edge was directed into *both* of its vertices (as the pattern output allows) when in the ALARM graph it is directed into only one. The results are summarized below:

	# of omitted undirected edges	#false undirected edges	#orientation errors
Trial 1	3	0	5
Trial 2	3	0	6
Trial 3	3	2	7

The implementation used did not determine the adjacency sets lying on undirected paths between two variables because in this case with correlation data it was computationally cheaper to determine the partial correlations for all subsets of A_{ab} than to keep track of $A_{ab} \cap U_{ab}$. With discrete count data for which the determination of conditional independence relations is more computationally demanding, the alternative procedure described in our statement of the algorithm might be faster. For example, for one pair of vertices in the network, A_{ab} consists of 8 vertices while $A_{ab} \cap U_{ab}$ consists of only two vertices.

The comparison of ten seconds for the PC algorithm with 22 and one half hours for the Kutato' algorithm should not be taken too literally, since the machines were quite different, and since without the assumption of linearity considerably more time would be required in numerical operations to determine conditional independence.⁵ Nonetheless, the PC algorithm appears to be very fast and reliable for sparse graphs. For similar data from a similarly connected graph with 100 variables, the present implementation should require less than two minutes.

References

I. Beinlich, H. Suermondt, R. Chavez, and G. Cooper. (1989) The ALARM Monitoring System, Technical Report KSL 88-84, Knowledge Systems Laboratory, Medical Computer Science, Stanford University.

⁵It may in fact be the case that for large samples and variable sets the errors introduced by assessing conditional independence through partial correlations or other aggregate measures are adequately repaid in time savings.

E. Herskovits and G. Cooper, Kutato': An Entropy-Driven System for Construction of Probabilistic Expert Systems from Databases, Proceedings of the Sixth Conference on Uncertainty in AI, Cambridge Mass, 1990.

P. Spirtes, C. Glymour and R. Scheines, 1990a, Causality from Probability, in G. McKee ed., *Evolving Knowledge in Natural and Artificial Intelligence*, Pitman; and 1990b *Proceedings of the Oak Ridge Conference on Advanced Computing for the Social Sciences*.

T. Verma and J. Pearl, 1990, Equivalence and Synthesis of Causal Models, Proceedings of the Sixth Conference on Uncertainty in AI, Cambridge.