

## **Automating Path Analysis for Building Causal Models from Data\***

Paul R. Cohen, Adam Carlson,  
Lisa Ballesteros, Robert St. Amant

Computer Science Technical Report 93-38

Experimental Knowledge Systems Laboratory  
Department of Computer Science  
University of Massachusetts  
Amherst, MA 01003

### **Abstract**

Path analysis is a generalization of multiple linear regression that builds models with causal interpretations. It is an *exploratory* or *discovery* procedure for finding causal structure in correlational data. Recently, we have applied statistical methods such as path analysis to the problem of building models of AI programs, which are generally complex and poorly understood. For example, we built by hand a path-analytic causal model of the behavior of the Phoenix planner. Path analysis has a huge search space, however. If one measures  $N$  parameters of a system, then one can build  $O(2^{N^2})$  causal models relating these parameters. For this reason, we have developed an algorithm that heuristically searches the space of causal models. This paper describes path analysis and the algorithm, and presents preliminary empirical results, including what we believe is the first example of a causal model of an AI system induced from performance data by another AI system.

---

\* This research was supported by DARPA-AFOSR contract F30602-91-C-0076.



regression with the data alone.  
effects may be used. Certain hypotheses may therefore be generated and  
experimented to see and the observation to manipulate variables to see  
natural consequences in an experiment. Nonexperimental means that the  
The term "nonexperimental" is because computing, because data are

quicker experimentally certain models. We must use such model  
value is a generalization of regression analysis that pro-  
duce certain models of our behavior. But analysis, how-  
ever, regression processes (e.g., regression analysis) to  
be found in regression models. Not only automatic  
like values but the certain, experimental knowledge was not  
before could affect the incidence of behavior and not  
certain relationships. For example, we know that the wind  
pressures increase, if blowing no experiment of their  
processes may be found. But experiment regression assumed  
certain and the: the wind speed or the number of times the  
times, such as "Which has more impact on the time to  
We intend to regression analysis to answer some ques-  
tioned experimentally with some the behavior of the process:  
the environment. However, we soon realize that we  
and we can't measure our behavior for our own and a re-  
sponse and time on its speed. At first, the process was  
behavior, but processes are not of data, and the process soon  
find the time: not time, but time, but time, but time,  
behavior now the other, which are semi-automatic, some  
processes and relationships. One thing, called the process,  
relationships, time and the activities of things such as  
works. The system, called process (Cohen et al., 80),  
experimentation of now a complex AI behavior system  
We developed the algorithm to help us discover certain  
value of correlation.

weaker situation than but analysis, which relies on evi-  
dence analysis: it relies on evidence of nondependence and a  
algorithm with a more general mathematical basis than  
relationships [23] may be used to develop a certain model  
times—essentially but analysis. But [21, 23] and  
LELAND system about the analysis of statistical data  
most similar to that of GILMORE et al. [81], who built the  
algorithm processes experimentally certain models. Our work is  
relationships of (behavioral) certain mathematical, our al-  
gorithms function finding algorithm processes function  
Leland et al., 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 101, 102, 103, 104, 105, 106, 107, 108, 109, 110, 111, 112, 113, 114, 115, 116, 117, 118, 119, 120, 121, 122, 123, 124, 125, 126, 127, 128, 129, 130, 131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 148, 149, 150, 151, 152, 153, 154, 155, 156, 157, 158, 159, 160, 161, 162, 163, 164, 165, 166, 167, 168, 169, 170, 171, 172, 173, 174, 175, 176, 177, 178, 179, 180, 181, 182, 183, 184, 185, 186, 187, 188, 189, 190, 191, 192, 193, 194, 195, 196, 197, 198, 199, 200, 201, 202, 203, 204, 205, 206, 207, 208, 209, 210, 211, 212, 213, 214, 215, 216, 217, 218, 219, 220, 221, 222, 223, 224, 225, 226, 227, 228, 229, 230, 231, 232, 233, 234, 235, 236, 237, 238, 239, 240, 241, 242, 243, 244, 245, 246, 247, 248, 249, 250, 251, 252, 253, 254, 255, 256, 257, 258, 259, 260, 261, 262, 263, 264, 265, 266, 267, 268, 269, 270, 271, 272, 273, 274, 275, 276, 277, 278, 279, 280, 281, 282, 283, 284, 285, 286, 287, 288, 289, 290, 291, 292, 293, 294, 295, 296, 297, 298, 299, 300, 301, 302, 303, 304, 305, 306, 307, 308, 309, 310, 311, 312, 313, 314, 315, 316, 317, 318, 319, 320, 321, 322, 323, 324, 325, 326, 327, 328, 329, 330, 331, 332, 333, 334, 335, 336, 337, 338, 339, 340, 341, 342, 343, 344, 345, 346, 347, 348, 349, 350, 351, 352, 353, 354, 355, 356, 357, 358, 359, 360, 361, 362, 363, 364, 365, 366, 367, 368, 369, 370, 371, 372, 373, 374, 375, 376, 377, 378, 379, 380, 381, 382, 383, 384, 385, 386, 387, 388, 389, 390, 391, 392, 393, 394, 395, 396, 397, 398, 399, 400, 401, 402, 403, 404, 405, 406, 407, 408, 409, 410, 411, 412, 413, 414, 415, 416, 417, 418, 419, 420, 421, 422, 423, 424, 425, 426, 427, 428, 429, 430, 431, 432, 433, 434, 435, 436, 437, 438, 439, 440, 441, 442, 443, 444, 445, 446, 447, 448, 449, 450, 451, 452, 453, 454, 455, 456, 457, 458, 459, 460, 461, 462, 463, 464, 465, 466, 467, 468, 469, 470, 471, 472, 473, 474, 475, 476, 477, 478, 479, 480, 481, 482, 483, 484, 485, 486, 487, 488, 489, 490, 491, 492, 493, 494, 495, 496, 497, 498, 499, 500, 501, 502, 503, 504, 505, 506, 507, 508, 509, 510, 511, 512, 513, 514, 515, 516, 517, 518, 519, 520, 521, 522, 523, 524, 525, 526, 527, 528, 529, 530, 531, 532, 533, 534, 535, 536, 537, 538, 539, 540, 541, 542, 543, 544, 545, 546, 547, 548, 549, 550, 551, 552, 553, 554, 555, 556, 557, 558, 559, 560, 561, 562, 563, 564, 565, 566, 567, 568, 569, 570, 571, 572, 573, 574, 575, 576, 577, 578, 579, 580, 581, 582, 583, 584, 585, 586, 587, 588, 589, 590, 591, 592, 593, 594, 595, 596, 597, 598, 599, 600, 601, 602, 603, 604, 605, 606, 607, 608, 609, 610, 611, 612, 613, 614, 615, 616, 617, 618, 619, 620, 621, 622, 623, 624, 625, 626, 627, 628, 629, 630, 631, 632, 633, 634, 635, 636, 637, 638, 639, 640, 641, 642, 643, 644, 645, 646, 647, 648, 649, 650, 651, 652, 653, 654, 655, 656, 657, 658, 659, 660, 661, 662, 663, 664, 665, 666, 667, 668, 669, 670, 671, 672, 673, 674, 675, 676, 677, 678, 679, 680, 681, 682, 683, 684, 685, 686, 687, 688, 689, 690, 691, 692, 693, 694, 695, 696, 697, 698, 699, 700, 701, 702, 703, 704, 705, 706, 707, 708, 709, 710, 711, 712, 713, 714, 715, 716, 717, 718, 719, 720, 721, 722, 723, 724, 725, 726, 727, 728, 729, 730, 731, 732, 733, 734, 735, 736, 737, 738, 739, 740, 741, 742, 743, 744, 745, 746, 747, 748, 749, 750, 751, 752, 753, 754, 755, 756, 757, 758, 759, 760, 761, 762, 763, 764, 765, 766, 767, 768, 769, 770, 771, 772, 773, 774, 775, 776, 777, 778, 779, 780, 781, 782, 783, 784, 785, 786, 787, 788, 789, 790, 791, 792, 793, 794, 795, 796, 797, 798, 799, 800, 801, 802, 803, 804, 805, 806, 807, 808, 809, 810, 811, 812, 813, 814, 815, 816, 817, 818, 819, 820, 821, 822, 823, 824, 825, 826, 827, 828, 829, 830, 831, 832, 833, 834, 835, 836, 837, 838, 839, 840, 841, 842, 843, 844, 845, 846, 847, 848, 849, 850, 851, 852, 853, 854, 855, 856, 857, 858, 859, 860, 861, 862, 863, 864, 865, 866, 867, 868, 869, 870, 871, 872, 873, 874, 875, 876, 877, 878, 879, 880, 881, 882, 883, 884, 885, 886, 887, 888, 889, 890, 891, 892, 893, 894, 895, 896, 897, 898, 899, 900, 901, 902, 903, 904, 905, 906, 907, 908, 909, 910, 911, 912, 913, 914, 915, 916, 917, 918, 919, 920, 921, 922, 923, 924, 925, 926, 927, 928, 929, 930, 931, 932, 933, 934, 935, 936, 937, 938, 939, 940, 941, 942, 943, 944, 945, 946, 947, 948, 949, 950, 951, 952, 953, 954, 955, 956, 957, 958, 959, 960, 961, 962, 963, 964, 965, 966, 967, 968, 969, 970, 971, 972, 973, 974, 975, 976, 977, 978, 979, 980, 981, 982, 983, 984, 985, 986, 987, 988, 989, 990, 991, 992, 993, 994, 995, 996, 997, 998, 999, 1000.

# 1. INTRODUCTION

then take sums:

$$\begin{aligned} \lambda X^3 &= \beta^1 X^1 X^3 + \beta^5 X^5 X^3 + \beta^3 X^3 \\ \lambda X^5 &= \beta^1 X^1 X^5 + \beta^5 X^5 X^5 + \beta^3 X^3 X^5 \\ \lambda X^1 &= \beta^1 X^1 X^1 + \beta^5 X^5 X^1 + \beta^3 X^3 X^1 \end{aligned}$$

producting three equations:

model producting by each of the terms on the right hand side  
of the regression equation. First, multiplying the prediction  
Multiplying linear regression finds the coefficients to con-  
strain in  $X^1$ .

was twice the influence on  $X$  as a standard deviation  
 $\beta^1 = 4$ ,  $\beta^5 = 8$ , then a standard deviation change in  $X^5$   
change in  $\lambda$ . Thus, the coefficients are combinations: if  
each deviation,  $\lambda$ , producting  $\beta^1$  standard deviation  
relation of this model is that a change in  $X^1$  of one stan-  
dard deviation are deviated with probability (e.g., the inter-  
prediction model:  $\lambda = \beta^1 X^1 + \beta^5 X^5 + \beta^3 X^3$  (standardized  
regression model with three predictor variables. Here is the  
how to solve for coefficients in a standardized linear re-  
gression. For convenience and simplicity, we will show  
from the variable its mean and deviation by its standard  
A standardized variable is what you get by subtracting  
models are often constructed for standardized variables,  
measured in miles. For this and other reasons, regression  
in fact, we will be a factor of 2580 (e.g., if  $x$  is  
of the predictor variables: for example, if  $x$  is measured  
regression coefficients depend on the units of measurement  
variable  $x$  producting  $\beta$  units change in  $\lambda$ . Thus, the re-  
relation of such a rule is that a unit change in a predictor  
value, units of the form  $\lambda = \beta^1 x^1 + \dots + \beta^k x^k + \alpha$ . The inter-  
relations and relationships for more than one predictor vari-  
Multiplying linear regression finds least-squares (e.g.,

$$\begin{aligned} \text{minimize } \sum (\lambda - \lambda^i)^2 \\ \text{single linear regression finds a line } \lambda = \beta x + \alpha \text{ that mini-} \\ \text{mizes of predicted values from actual values. That is,} \\ \text{square line is one that minimizes the sum of squared de-} \\ \text{viations of } x \text{ to a best-fitting line } \lambda. \text{ A least-} \\ \text{square line finds a least-squares line relating a single predic-} \\ \text{tor } x \text{ to } \lambda \text{ with regression. Single linear re-} \\ \text{gression analysis is a generalization of multiple linear re-} \end{aligned}$$

# 2. BACKGROUND: REGRESSION

ing it to statistical data.  
which the behavior of the algorithm was proved by applying  
algorithm to Proenix data, and a statistical experiment in  
two experiments, an informal one in which we applied the  
regression 2 describes our algorithm. Regression 3 describes  
regression 4 where we illustrate a but analysis of Proenix  
data to regression 3, where we introduce but analysis, or  
behavior who are similar with regression analysis which  
certain experiments of now the Proenix process work:  
we describe below, we may be able to discover other  
of Proenix by hand, and by automatically but analysis as

not worth the effort. Path analysis would be identical to linear regression and the path model in Figure 1. If this were all we could do, causal interpretation of the normal equations in terms of model and least-squares fit to our data. We also have a point, and to beta coefficients that make the prediction  $\lambda = \beta^1 X^1 + \beta^2 X^2 + \beta^3 X^3$  gives rise to a set of normal equations that we have described how a prediction model

correlated with  $X^1$ : we will not consider this case here.) might be due to an unmeasured or latent variable that is fixed on  $X$  is attributable to  $X^1$ . (Alternatively, the effects are fixed and  $X^1$  is varied; in such an experiment, the effect of the control group is in an experiment in which  $X^2$  and  $X^3$  this sense, beta coefficients provide a statistical version of path to  $X$  when only  $X^1$  is systematically varied. In direct variables are fixed. You can interpret  $\beta^1$  as what effect of a predictor variable on  $X$  when all the other predictor variables are fixed. You can interpret  $\beta^2$  as what effect of a predictor variable on  $X$  when all the other predictor variables are fixed. The causal interpretation of causal effects without allowing regression correlations and direct causal relationships in Figure 1. By convention, links of the path. The second and third normal equations (either correlations or betas) along the constituent where the weight of a path is the product of the coefficients given by the sum of the weights of three paths in Figure 1, indirect path through  $X^3$ . Thus, the correlation  $\lambda X^1$  is through  $X^2$  to  $X$ ; and the third term is represented by the second term is represented by the indirect path from  $X^1$  in Figure 1 by the direct path between  $X^1$  and  $X$ : the  $\lambda X^1 = \beta^1 + \beta^2 \lambda X^2 + \beta^3 \lambda X^3$ . The  $\beta^1$  term is represented in Figure 1. Consider the first normal equation and they have an interesting interpretation, illustrated in The three equations (1) are called the normal equations,

this demonstration in [1, 2] or any good statistics text. is a least-squares rule, but the interested reader can find these coefficients guarantee that  $\lambda = \beta^1 X^1 + \beta^2 X^2 + \beta^3 X^3$  three unknown beta coefficients. We have not shown that clearly, with these three equations we can solve for the

$$\begin{aligned} \lambda X^3 &= \beta^1 \lambda X^2 + \beta^2 \lambda X^3 + \beta^3 \\ \lambda X^2 &= \beta^1 \lambda X^2 + \beta^2 + \beta^3 \lambda X^3 \\ \lambda X^1 &= \beta^1 + \beta^2 \lambda X^2 + \beta^3 \lambda X^3 \end{aligned} \tag{1}$$

be written in terms of correlations: (divided by  $\lambda$ ) is the variance, the previous equations can ignore here) and the square of a standardized variable variables is their correlation (divided by  $\lambda$ , which we can now, because the sum of the product of two standardized

$$\begin{aligned} \sum \lambda X^3 &= \beta^1 \sum X^1 X^3 + \beta^2 \sum X^2 X^3 + \beta^3 \sum X^3 \\ \sum \lambda X^2 &= \beta^1 \sum X^1 X^2 + \beta^2 \sum X^2 + \beta^3 \sum X^3 X^2 \\ \sum \lambda X^1 &= \beta^1 \sum X^1 + \beta^2 \sum X^2 X^1 + \beta^3 \sum X^3 X^1 \end{aligned}$$

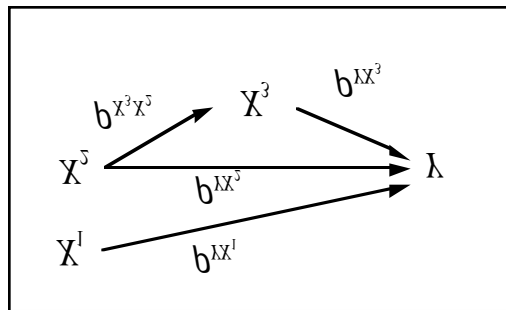
then another variable  $X$  (as  $X^1, X^2, X^3$  point to  $X$  in Fig. 2) beta coefficient. If several variables  $A, B, \dots$  all point to  $X$ . The second step is to fit each one as a correlation or a standardizing path diagram. Imagine Figure 2 is the result. This step is to propose a prediction model and a correlation prediction model is a least-squares fit to the data. The of why these steps lead beta coefficients that ensure that path analysis. See [Cohen] for a more formal discussion. What follows is an informal description of the steps in data.

that make the prediction model a least-squares fit to the data, and if the latter, path analysis solves for the values tells us whether these coefficients are correlations or related with beta coefficients, denoted  $\beta$ . Path analysis model is  $\lambda = \beta \lambda X^1 + \beta \lambda X^2 + \beta \lambda X^3$ . The arcs are  $X^1, X^2$  and  $X^3$  exist. The corresponding prediction

tion of  $X$  on  $X^1, X^2$  and  $X^3$ .

ground to a multiple linear regression-

Figure 2: A path model that does not corre-



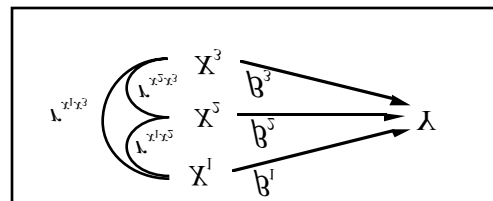
ences  $X$  through  $X^3$ , and no intercorrelations between  $X^1$  and  $X^2$  directly cause  $X$ , and  $X^2$  also indirectly influence  $X$ . For example, Figure 2 shows a model in which coefficients that ensure the model is a least-squares fit to our data. The power of path analysis is that we can specify virtually

### 3. PATH ANALYSIS

on  $X^1, X^2$  and  $X^3$ .

the multiple linear regression of  $X$

Figure 1: The path model that corresponds to





ply based on imprecise measure of goodness was the error term. In fact, when we constructed the Proemix model  $\chi^2$  than a highly-connected model, but will often be better than a model with necessarily large a lower influence are isolated, not distributed around a network of causal relations among variables as possible, so causal network causal story, and that includes as few correlations and goodness. The researcher wants a model that is a binary not necessarily the researcher's binary measure of goodness of that is not. So  $\chi^2$ , the statistical measure of goodness, is not, and so, broadly influence the amount of influence we know that will spread influence the rate at which the  $\chi^2$  measure is "causal process" of  $\chi^2$ . In fact, bounding directly to  $\chi^2$ , but we know that and  $\chi^2$  as causal at the same "level", that is, both regression analysis of the Proemix data sets  $\chi^2$  causal model of any system we analyze. For example, a Bayesian model. It is if likely to be a Bayesian variable: every variable directly influences  $\chi^2$ . It is not a (e.g., Fig. 1). No wonder this model accounts for so much correlation and all but directly to the dependent variable regression model, in which all independent variables are no model accounts for more of the variance in  $\chi^2$  than the complete. However, for any set of independent variables, influence  $\chi^2$ , and our understanding of  $\chi^2$  is therefore in  $\chi^2$  is low, then other variables besides  $\chi^2, \chi^2, \dots$  in  $\chi^2$  accounted for by the independent variables  $\chi^2, \chi^2, \dots$ .  $\chi^2$ , the measure of variance in the dependent variable, is "good." A common statistical measure of goodness is the mean square deviation two senses in which a model can

and this algorithm is constructing the evaluation function. The most complex issue to be addressed in hybrid model or we can make no more significant improvement. We continue until either we have reached an acceptable model, evaluate the new model, and add it to the we select the best modification in the list, hybrid if to its modifications to those models. After that in the search space we maintain a list of all the models and all possible combinations of just the dependent variables. During the search (possibly new) variable. We begin with a graph the addition of an arc to one variable in the model from both model constructs a state, while the sole operator is. The algorithm uses a form of best-first search. A single search space is of size  $2^N$ .

where  $n$  is the number of variables in the model. Thus the hybrid model then as an adjacency matrix of size  $N \times N$ , constraints described in section 3. We can represent any hybrid models. Hybrid models are graphs which satisfy the. The search space for the algorithm is the set of all possible combinations such an algorithm.

automatically by a heuristic search algorithm. This search whether models like the one in Figure 3 can be generated. The question that motivated the current research is

**MODELS**  
**2. AUTOMATIC GENERATION OF BATH**

evaluation of the resulting model. It is to be noted, the predicted correlation error and the weighted sum of the actual correlation between the variables. In the current implementation we use a. One evaluation of a modification is a function of these

- heuristic is appropriate for the problem.
- a Bayesian sense of values for a domain independent
- Bayesian variables are likely/unlikely causes of other
- Bayesian variables are independent.

forms of knowledge. These include knowledge that the third class represents domain/brother dependent

- Various adjustments to  $\chi^2$ .
- bandwidth:
- the "attenuation" of the variables and arcs in the

model are. Others which we have considered, but not yet imple-

- the total number of variables and arcs.
- a new link.
- the correlation between the variables being connected (e.g., the ratio of arcs which don't introduce a new variable).
- Bayesian (i.e. the ratio of variables to arcs) or the and the predicted correlation matrix for the model): adjusted error between the actual correlation matrix
- the predicted correlation error (minimize the total
- $\chi^2$ , the statistical measure of goodness.

problem. Some heuristics we have tried are weighting of the heuristic may depend on the Bayesian. These apply to any hybrid analysis problem; however, the. The second class contains domain independent heuristics:

- a dependent variable may have no outgoing arcs.
- from every variable to a dependent variable.
- there must be a path (in the sense of Wright's rules)
- no cycles are allowed.

that what we might call genetic criteria: heuristics into three general classes. The first class good models and to evaluate a hybrid model. We can group. We rely on a variety of heuristics to guide search toward evaluation function.

evaluation function as a term in the modification neighborhood. This is achieved by including the model evaluation heuristics have brought the search into the right function should dominate the search once the modification. These are the modification heuristics. A model evaluation move the search into that general region of the space. Considered together in model space, a few heuristics can and model evaluation. Assuming that good models are. Another distinction is between modification evaluation never even calculated  $\chi^2$ . between estimated and empirical correlations, and we

speed automatically.

changes model of a complex software system over generations, if that process what we believe is the first algorithm the algorithm can slowly, and in some ways produce our model in Figure 3.

should not be too surprised that the algorithm did not converge, and not concerned at all about  $R_5$ . Thus we mainly about the errors between estimated and empirical we, when we run Figure 3 by hand, were concerned by and modification scoring functions that value  $R_5$ , while however, we note that the algorithm used model scoring function on Figure. In defense of the algorithm, since of  $R_{LK}$  on function and the influence of the models is that neither recognizes the important information from is not argued. A disagreement aspect of this in our experiment, mindless and  $R_{LK}$  clearly so brooding this curiosity we realized that due to a sampling dependent variables set by the experimenter. However, mindless cases  $R_{LK}$  when, in fact, these are independent. A set aspect of the models is that both set factors, which are measures of behavior of the process time: the mind speed is causally process of the other causal order of mindless, function, and Figure respects. A good aspect of the models is that they get the mind to comment them, but they are also flawed in some ways were the ones shown in Figure 4. These models have judgments Exhibit II+ List Machine, its best two algorithm terminated, after four hours work on a Texas to limit the solution space to 820 models. When the algorithm-scoring functions were sufficiently better than  $\Sigma_{R_5} - \Sigma_{R_1} = 1^{\circ}048^{\circ}210^{\circ}$ , but the model-scoring and the dependent variable. The search space of models was function from the data set). We designed function function function, and function (we dropped the process parameter, specifically, mindless  $R_{LK}$ . We provided the algorithm with data from 112 trials of

6. EXPERIMENT 1

this we used artificial data. The second experiment was more expansive: for the algorithm could generate a model from the process performance. In the first experiment we tested whether performance of our algorithm and probe factors that affect. We have run several experiments to demonstrate the per-

6. EXPERIMENT 2

increase efficiency. changes to each model's evaluation. This should greatly brooding the algorithm by incrementally adjusting dominant search cost, we are currently working on increasing newly generated model. Because these calculations algorithm we calculate these parameters from search for correlation estimates between variables. In the basic form, including both coefficients in the model and the computation required for the heuristic evaluation function. The modification evaluation step takes a good deal of

the time from each independent variable to the dependent measured both score. For each rule model, make a list of all (roles) by two criteria:

We evaluated the performance of our algorithm (step 2, checked in the Appendix.

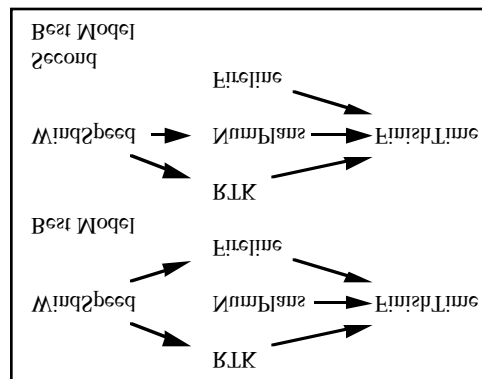
model specified in step 5. The details of the process are state data with the same correlational structure as the both joining this brooding (step 3, roles) ensures that we generate correlational value function, specified value, function: generate a column of numbers so that all their specified value, say, '8'. Now make the brooding more such that the correlation of the samples is a function, from a normal (Gaussian) distribution with mean zero one sets you to generate two samples of numbers drawn steps 3 and 2 require some explanation. Imagine some-

- 2. Determine how well our algorithm discovered the models if proposed.
- 4. Summarize the data to our algorithm and record the
- 3. Generate data consistent with the weights in step 5, resulting model should exceed 'd'.
- 5. Randomly assign weights to the links in the both produce a both model.
- 1. Randomly fill some cells in an adjacency matrix to

followed these steps: brooding could discover the rule models. Specifically, we that represented "rule" and tested how frequently our algorithm an experiment in which we constructed both models "rule" model. To address these and other questions we How do we know whether the algorithm is finding the source depend on the number of data for each variable? on the sample variance for each variable? Does performance: Does the performance of the algorithm depend. The first experiment raised more questions than it answered.

6. EXPERIMENT 3

models found by the algorithm. Figure 4: The best and second best process

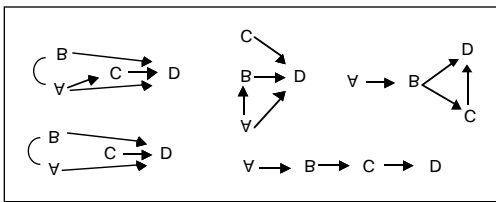


link weights to the true model (step 2, above), the model is lower than expected. This difference, after adjusting for it we fit the best model to those data, the resulting  $R^2$  in fact, generates data that have the expected correlations of each best model (as described in the Appendix), we do when we generate data to match the correlation structure. The reason for the disparity is actually quite simple: found the true model, which says about our expected score and the true model score was less in the true best model. However, the disparity between the best model criteria led the algorithm suggest, away from the true model. Our immediate inclination was to say the scoring the algorithm was better than the average score of the true in all trials, the average score of the best model found by

our criteria for terminating a search can be improved. model is very rare but too soon. This in turn suggests that which suggests that one reason they failed to find the true

model.

Figure 2: Five best models used to test the al-



gorithm. We would like to place more stock in two measures of the proportion of models that were possible. This that explored 134,000 models, on average. This is about four trials. When the algorithm did find the true model, it the algorithm did not find the true model, are presented in table 2. Summary statistics for these trials, and trials in which the algorithm found the true model in 100 of the 540 it-

algorithm ran 540 times. and in each of the 4 conditions just described. Thus, the produced 60 models. We ran our algorithm on these vari-corrrelation structure established in step 2, above. This of data for variables A, B, C and D that conform to the For each model we generated 15 variants, that is, 15 sets true best models with four variables, shown in Figure 2, in the samples would be expected. We generated five sizes became small, because the effects of outlier values mass of the algorithm might deteriorate when sample 20, and 100 data per sample. We thought that the best-berment we looked at four levels of the factor: 10, 30, data included 512 values of each variable, but for this ex-ize on the performance of the algorithm. The Phoenix Our experiment was designed to test the effects of sample exist in the model discovered by our algorithm; the links between variables; what fraction of these links shared link score; For each true model, make a list of all discovered by our algorithm; variable; what fraction of these best exist in the model

the best found model, or the shared best or shared link number of models expanded by the algorithm, the score of shared; sample size affects to have no impact on the best variables were the true measures in Table 2, of the number of data points in our samples. The degen-We ran a one-way analysis of variance to find the effects given by the shared best score and shared link score.

proved, but it performs better than the true impression two links on average. Clearly, the algorithm can be im-found model differed from the true one by little less than shared link scores are low, they suggest that the best score would be 25%. So although the shared best and model was just two links, when the average shared link the disparity between the true model and the best found model. These numbers are not high. On the other hand, if the links, in the true model are also in the best found was. On average, 20-21% of the best, and 20-24% of in which the true model was found from those in which it Unfortunately, neither score strongly differentiates trials model that are also in the best found model. the fraction of the links connecting variables in the true be 0.6. Another criterion is the shared link score, which is from A, B and C to D, then the shared best score would was, say, the top-right model, which has only three bests from A, B and C to D, so if the best found model model. For example, the top-left model in Figure 2 has 2 bests in the best found model that also exist in the true found model. The shared best score is the proportion of structural similarity between the true model and the best We would like to place more stock in two measures of

Shared link score	.245	.163	.201	.144
Shared best score	.015	.505	.262	.133
True model score	.841	.528	.808	.741
Best model score	.674	.551	.624	.618
Number of models	134,000	63,300	80,000	57,700
	Mean	Std.	Mean	Std.
	model found the true Trials that did		true model not find the Trials that did	

Table 2: Summary statistics from 540 trials.

model. claim that our algorithm found or failed to find the "best" score. This problem will have to be solved before we can we generated the data is not the model with the highest Consequently, in most of our trials, the model from which an artifact of our procedure for generating data. accounts for, say, 80% of the variance in the data. This is and fit the data to the model, we find that the model only 5% might be, say, .02, but after we generate data (step 3)



done to solve both of these problems:
   
bays and sparse jinks, were not ideal. More needs to be
   
row group we measure closeness? Our criteria, sparse
   
not find exactly the "correct" model, but one way like it
   
ing function. Second, in the event that the algorithm does
   
ry the statistical combination (i.e.,  $V_5$ ) of the model scor-
   
which we generated the data, might be given a low score
   
in this case, the "correct" model, that is, the one from
   
better explained by a model other than the original model.
   
First, it often happens that we generated data that were
   
of the experiment is some, two technical problems arose:
   
algorithm could find the original model. While the logic
   
correlation structure of the model, then see whether our
   
tion a bay model, then generate data consistent with the
   
usually straightforward. Our procedure was to first decide
   
generated the data—then evaluation of goodness is tech-
   
model—the model that corresponds to the given that
   
it is good we mean that the algorithm finds the "correct"

All of these kinds of models are good by some criteria
   
and dependent variables, such as the models in Figure 4,
   
produce causal models with nodes between the predictor
   
relationship are weighted heavily, then the algorithm will
   
it errors between observed and empirical correlations and
   
tends the intercorrelations among the predictor variables,
   
keep the fit causal structure of regression models, but
   
model is also weighted heavily, the algorithm tends to
   
dependent variables and are fully intercorrelated. If vari-
   
tion models, in which all variables both directly to the
   
derived, then the algorithm will always generate regres-
   
might expect if  $V_5$  is weighted heavily and relationship is
   
us to adjust the relative importance of the criteria. As you
   
The algorithm is currently biased towards, which variables
   
be evaluated by both statistical and biological criteria.
   
produce good models. As we noted earlier, goodness can
   
More important in the near term is whether the algorithm
   
models in a reasonable time:

can improve the algorithm enough to find large causal
   
number of variables. It remains to be seen whether we
   
of the search space is exponential in the square of the
   
number significant to the run time. Still, the complexity
   
generated, which takes a lot of space, so finding com-
   
bined: for example, we keep all models that have been
   
works well. The algorithm's speed can certainly be im-
   
let to establish criteria for demonstrating clearly that it
   
models from data. If works slowly, however, and we have
   
sense the algorithm "works," that is, if generates causal
   
results from experiments with the algorithm. In a crude
   
We have described an algorithm for bay analysis and that

## 7. CONCLUSION

few data:
   
shows that the algorithm can work even with relatively
   
sample size of the method for generating data. Thus, it
   
be a statistical artifact, due entirely to the sensitivity to
   
on the score of the true model. This effect turned out to
   
scores. But the sample size has a large significant effect

size and forming  $\underline{x}$  from them. Since each  $x^b$  is sampled
   
values generating  $x^1, x^2, \dots, x^b$  independent normal vari-
   
The method for generating  $\underline{x}$  from this distribution in-

$$\underline{x} = A\underline{z} + \underline{\mu} \tag{5}$$

that
   
formation of  $b$  independent normal variables in  $\underline{z}$ , such
   
The distribution of  $\underline{x}$  can be represented as a linear trans-

$$\chi(x) = (\sigma A)^{-1} \prod_{i=1}^b \exp\{-\frac{1}{2}(\underline{x} - \underline{\mu})^T A^{-1}(\underline{x} - \underline{\mu})\} \tag{1}$$

$A$  (non-singular) is defined by the b by
   
matrix having mean  $\underline{\mu}$  and variance-covariance matrix
   
multivariate normal distribution. The  $b$ -variate normal dis-
   
tributed linear combination of the  $b$  components of  $\underline{x}$  has a
   
multivariate normal distribution if and only if every non-
   
 $A$   $b$ -dimensional random vector  $\underline{x}$  is defined to have a
   
matrix used to calculate that variate.

relationship structure underlying the variance-covariance
   
efficiency and because it returns a variate having the cor-
   
tion a multivariate normal distribution was chosen for its
   
the model. The procedure for the generation of variables
   
erated for a model have the same correlational structure as
   
relation matrix is calculated. It is essential that data gen-
   
After initialization of the bay weights, the predicted cor-

- $\delta < V_5 \leq 1.0$   
  mean one
- all predicted correlations for the model must be less
- $\rho \geq$  correlational bay weights  $< 1$
- $\rho \geq$  causal bay weights  $< 1$

are randomly generated under the following constraints:
   
pose of data generation, the bay weights of a given model
   
describes the relationships among its variables. For the per-
   
lation matrix, calculated from its bay weights, which de-
   
Every bay model has associated with it a predicted corre-

## APPENDIX: DATA GENERATION

ods:
   
who seek to understand AI systems with statistical meth-
   
ing system, and promises to be a valuable tool for those
   
gorithm generated causal models of a complex AI sys-
   
ever difficult it is to evaluate, the fact remains that the al-
   
causal models of systems. However, however, it is, how-
   
promise the algorithm holds for automatically finding
   
Still, we are encouraged by our results, and by the
   
clearly.

models. Clearly, we would prefer this curve to rise less
   
erage, but with five variables it searched thousands of
   
ables, the algorithm searched roughly 100 models on av-
   
increases with the number of variables. With four vari-
   
cursive. We also need to know how the solution space
   
ble data affects performance, but the results were incon-
   
an experiment to see whether the variance of our sam-
   
affect the performance of the algorithm. For example, we
   
We have yet to answer many questions about factors that

a  $p \times n$  data matrix having  $n$  data points for each variable.  $\mathbf{y}$  is repeated  $n$  times, each vector  $\mathbf{x}_i$  forming the columns of  $\mathbf{X}$ .  $\mathbf{y}$  is substituted into equation (5) to calculate  $\hat{\mathbf{x}}$ . The process and  $\hat{\mathbf{x}}$  is chosen to be a zero vector. These values are matrix  $\mathbf{A}$  is formed from a Choleski decomposition of  $\mathbf{V}$ . The variance-covariance matrix is equal to the correlation from the normal distribution with standard deviation 1.

Causation, Prediction and Search. Zbigniew-Verlag.  
Zbigniew B., Gilmore, C. and Zschaler, B. (1993)

The Mill Press. 858-833.  
Eighth National Conference on Artificial Intelligence:  
the function-learning algorithm. In Proceedings of the  
Zschaler, C. 1990. A broken domain-independent scien-

Er Gauderthal, EG. 141-120.  
of the Fourth International Workshop on AI and Statistics:  
graphs admit a causal interpretation; Preliminary Papers  
Bert, J. & Wernup, N. 1993. When can association

Morgan Kaufman. 441-425.  
Conference. J. Allen, K. Elkes, & E. Sandewall (Eds):  
Reasoning: Proceedings of the Second International  
Workshop Principles of Knowledge Representation and  
Bert, J. & Wernup, J. 1991. A theory of inferred causal-

PLCC. 1992. Path Analysis-A Primer. Boxwood Press.

Explorations of the Causal Process. The Mill Press.  
J.M. 1981. Scientific Discovery: Computational  
Langley, B., Simon, H.A., Bradshaw, G.T. & Zytkow,  
Planning Systems. Morgan Kaufman. 109-112.

Proceedings of the First International Conference on AI  
and Success and Task Duration in the Phoenix Planner.  
Hart, D.M. & Cohen, B.B. 1985. Predicting and explain-  
Discovering Causal Structure. Academic Press.

Gilmore, C., Zschaler, B., Zbigniew, B. & Kelly, K. 1981.  
tern. Machine Learning 1(1): 391-401.  
dynamizable and undynamizable discovery: the ABCD2 sys-  
Falkenhainer, B.C. & Michalski, R.S. 1989. Integrating

10(3): 35-48.  
ments for agents in complex environments. AI Magazine,  
1988. That's all fine: Understanding the design routine.  
Cohen, B.B., Greenberg, M.T., Hart, D.M. & Howe, A.E.

188.  
Workshop on AI and Statistics. Er Gauderthal, EG. 182-  
appear in Proceedings of the Fourth International  
an autonomous agent in a complex environment. In  
Cohen, B.B. & Hart, D.M. 1993. Path analysis models of

Forecasting.  
Cohen, B.B. Empirical Methods for Artificial Intelligence.  
Newbury Park, CA.

Asker, H.B. 1983. Causal Modeling. Sage Publications.

**REFERENCES**

person-  
mental hypotheses notwithstanding any copyright notice  
intended to reproduce and distribute reprints for govern-  
E30905-91-C-0079. The United States Government is au-  
This research was supported by DARPA-AFOSR contract  
development of these ideas.

We would like to thank Glenn Shafer for his help with the

**Acknowledgments**

The Mill Press. 888-884.  
the Eighth National Conference on Artificial Intelligence:  
discovery in a chemically laboratory. In Proceedings of  
Zytkow, J.M., Shim, J. & Hussain, A. 1990. Automated

