A Woefully Incomplete History of Empirical Applications of Algorithmic Causal Discovery (ACD)

#### UAI 2023

#### Workshop:

The History and Development of Search Methods for Causal Structure

Richard Scheines Carnegie Mellon University

## Literature searches

Thanks to:



• 2013 Workshop at CMU on empirical applications of ACD:

https://www.cmu.edu/dietrich/philosophy/events/workshops-conferences/causal-discovery/index.html

• Work since by authors from Workshop at CMU:

https://docs.google.com/document/d/1nv84oQYy9YAXvn5RxNeGcaGwsmv\_RILKE50hHNWWxcw/edit

• Other Published Applications of "Causal Discovery" Methods 2003 – 2023:

https://docs.google.com/document/d/1mcjLIFk8JRUVN99rfv3BGAPUSoCGtC13ZiY1pY4wUMk/edit?usp=sharing

## The History of Applying ACD to Empirical Data

- ACD tools ~ 40 years : "Structure Learning", "Causal Discovery", etc.
- Now: dozens of algorithms, toolkits, full suites available, in Java, Python, R, etc.
- Massive # of algorithm + simulation studies / small # empirical applications why?
- ACD requires assumptions, e.g., about:
  - Time
  - Feedback
  - Distributions
  - Parametric families
  - Sampling
  - Latent confounders
  - Measurement / measurement error

## Still – there are 100s of applications –

How to categorize them for this presentation?

- Algorithm
- Assumptions
- Discipline
- Experimental follow-up
- Scientific purpose

# Common Scientific Purposes of Using of ACD on Empirical Data

- 1. Finding alternatives .... to:
  - 1. Theoretically specified causal hypotheses
  - 2. Regression based causal inferences /

standard factor analysis latent variable models

2. Generating candidate causes

3. Mechanism Discovery

Presentation goal: see the rough outlines of the landscape,

So: simple and shallow and fast

Finding alternatives .... to:

1. "Theoretically" specified causal hypotheses

## Charitable Giving

What influences giving? Specificity? Sympathy? Impact?

"The Donor is in the Details", (2013) Organizational Behavior and Human Decision Processes, Issue 1, 15-23, C. Cryder, with G. Loewenstein, R. Scheines.

N = 94

TangibilityCondition	[1,0]	Randomly assigned experimental condition
Imaginability	[17]	How concrete is the scenario
Sympathy	[17]	How much sympathy for target
Impact	[17]	How much impact will my donation have
AmountDonated	[05]	How much actually donated

## **Theoretical Hypothesis**



### **Estimated Model**



### **Search Alternative**



## Common Scientific Purposes of Using of ACD on Empirical Data

- 1. Finding alternatives .... to:
  - 1. Theoretically specified causal hypotheses
  - 2. Regression based causal inferences /

standard factor analysis latent variable models

**Regression for Causal Inference** 

- X a prima facie cause of Y? (i) plausible, and (ii) X and Y associated
- Regress Y on X and potential confounders Z<sub>i</sub>:
  - Z<sub>i</sub> associated with X, and
  - Z<sub>i</sub> associated with Y, and
  - Z<sub>i</sub> prior to X

#### Regression is bad for causal inference – ACD is better

No practicing social scientist or statistician seems to know this – or take it seriously



#### Prima facie cause? Yes



 $X \rightarrow Y$  ??

## Confounders to control for in a regression?

## Z1 and Z2



 $X \rightarrow Y$  ??

## $X \rightarrow Y$ ?? Regression: Yes - *incorrect*!!



## $X \rightarrow Y$ ?? ACD: No - correct!!



## Foreign Investment

Does Foreign Investment in 3<sup>rd</sup> World Countries cause Repression?

Timberlake, M. and Williams, K. (1984). Dependence, political exclusion, and government repression: Some cross-national evidence. *American Sociological Review* 49, 141-146.

#### N = 72

- po degree of political exclusivity
- cv lack of civil liberties
- en energy consumption per capita (economic development)
- fi level of foreign investment

#### Foreign Investment

#### Correlations



#### Prima Facie:

Foreign Investment associated with LESS Political exclusion (more Political repression) Foreign Investment

### **Regression Results**

SE	(.058)	(.059)	(.060)
t	3.941	-2.99	14.6

Controlling for Economic Development (en) and lack of civil liberties (cv) flips the sign:

Foreign Investment causes More Political exclusion (LESS Political repression)

#### Case Study: Foreign Investment Alternative Models



There is no linear model with testable constraints (df > 0) in which FI has a positive effect on PO, that is not rejected by the data.



## Lead $\rightarrow$ IQ 1979 NEJM Study

- 2500 children's teeth collected and measured for lead exposure
- All children rated behaviorally by teacher



## 1979 NEJM Study

- 39 Potential confounding variables measured:
  - Socioeconomic status
  - Parental IQ
  - Mother's age at birth
  - Father's age at birth
  - Mother's level of education
  - Number of live birth's before sampled child
  - Etc.

## Lead and IQ: Variable Selection



Final Variables (Needleman)

- -lead baby teeth
- -fab father's age
- -mab mother's age
- -nlb number of live births
- -med mother's education
- -piq parent's IQ
- -ciq child's IQ

## Needleman Regression

- standardized coefficient
- (t-ratios in parentheses)
- p-value for significance

ciq =	143 lea	ad204 f	ab159 n	ılb + .219 m	ned + .237 r	mab + .247	piq
	(2.32)	(1.79)	(2.30)	(3.08)	(1.97)	(3.87)	
	0.02	0.09	0.02	<0.01	0.05	< 0.01	

All variables significant at .1  $R^2 = .271$ 

## **TETRAD** Variable Selection

#### Tetrad

mab \_||\_ ciq fab \_||\_ ciq

nlb\_||\_ciq | med



## Regressions

- standardized coefficient
- (t-ratios in parentheses)
- p-value for significance



TETRAD (
$$R^2 = .243$$
)  
ciq =  $-.177$  lead  $+.251$  med  $+.253$  piq  
(2.89) (3.50) (3.59)  
<0.01 <0.01 <0.01

# Common Scientific Purposes of Using of ACD on Empirical Data

#### 1. Finding alternatives .... to:

- 1. Theoretically specified causal hypotheses
- 2. Regression based causal inferences /

standard factor analysis latent variable models

#### Case Study: Stress, Depression, and Religion

MSW Students (N = 127) 61 - item survey (Likert Scale)

- Stress: St<sub>1</sub> St<sub>21</sub>
- Depression: D<sub>1</sub> D<sub>20</sub>
- Religious Coping: C<sub>1</sub> C<sub>20</sub>



#### Case Study: Stress, Depression, and Religion



### Case Study: Stress, Depression, and Religion



# Common Scientific Purposes of Using of ACD on Empirical Data

- 1. Finding alternatives .... to:
  - 1. Theoretically specified causal hypotheses
  - 2. Regression based causal inferences /

standard factor analysis latent variable models

#### 2. Generating candidate causes

3. Mechanism Discovery

### Candidate Gene Regulators

Which genes regulate flowering time in *Arabidopsis thaliana*?



\* Stekhoven DJ, et al. Causal stability ranking. *Bioinformatics* 28 (2012) 2819-2823.

## **Observational Data**

 n = 47 Arabidopsis thaliana gene expression profiles of 4-day old seedlings for which subsequent flowering time was also measured

• Affymetrix ATH1 arrays with expression measurements on 21,440 *A. thaliana* genes

### **Candidate Gene Selection**



#### Candidate Regulators of Flowering Time

- Considered the 25 genes that were ranked most likely to be causes of flowering time, according to the causal network analysis
- 5 of those 25 genes were known regulators of flowering
- 13 of those 25 genes were not known regulators and mutant seeds for each of them were available


Greenhouse experiments on flowering time



# **Experimental Details**

- Seeds were planted, and flowering time was measured in days to bolting
- Seeds types yielding 4 or more plants were considered viable for analysis (9)
- 4 of 9 had a statistically significantly shorter mean flowering time (p < 0.05) than the control, wild-type plants
- Correlational and other analyses to identify candidate genes were successful roughly at chance rates, i.e., they were worthless

# The Tumor-specific Driver Identification (TDI) Algorithm



> PLoS Comput Biol. 2019 Jul 5;15(7):e1007088. doi: 10.1371/journal.pcbi.1007088. eCollection 2019 Jul.

#### Systematic discovery of the functional impact of somatic genome alterations in individual tumors through tumor-specific causal inference

Chunhui Cai <sup>1</sup> <sup>2</sup>, Gregory F Cooper <sup>1</sup> <sup>2</sup>, Kevin N Lu <sup>1</sup> <sup>2</sup>, Xiaojun Ma <sup>1</sup>, Shuping Xu <sup>3</sup>, Zhenlong Zhao <sup>3</sup>, Xueer Chen <sup>1</sup> <sup>2</sup>, Yifan Xue <sup>1</sup> <sup>2</sup>, Adrian V Lee <sup>2</sup> <sup>3</sup> <sup>4</sup> <sup>5</sup>, Nathan Clark <sup>2</sup> <sup>6</sup>, Vicky Chen <sup>1</sup> <sup>2</sup>, Songjian Lu <sup>1</sup> <sup>2</sup>, Lujia Chen <sup>1</sup> <sup>2</sup>, Liyue Yu <sup>1</sup> <sup>2</sup>, Harry S Hochheiser <sup>1</sup> <sup>2</sup>, Xia Jiang <sup>1</sup> <sup>2</sup>, Q Jane Wang <sup>3</sup>, Xinghua Lu <sup>1</sup> <sup>2</sup> <sup>5</sup>

Affiliations + expand PMID: 31276486 PMCID: PMC6650088 DOI: 10.1371/journal.pcbi.1007088 Free PMC article

#### Abstract

Cancer is mainly caused by somatic genome alterations (SGAs). Precision oncology involves identifying and targeting tumor-specific aberrations resulting from causative SGAs. We developed a novel tumor-specific computational framework that finds the likely causative SGAs in an individual tumor and estimates their impact on oncogenic processes, which suggests the disease mechanisms that are acting in that tumor. This information can be used to guide precision oncology. We report a tumor-specific causal inference (TCI) framework, which estimates causative SGAs by modeling causal relationships between SGAs and molecular phenotypes (e.g., transcriptomic, proteomic, or metabolomic changes) within an individual tumor. We applied the TCI algorithm to tumors from The Cancer Genome Atlas (TCGA) and estimated for each tumor the SGAs that causally regulate the differentially expressed genes (DEGs) in that tumor. Overall, TCI identified 634 SGAs that are predicted to cause cancer-related DEGs in a significant number of tumors, including most of the previously known drivers and many novel candidate cancer drivers. The inferred causal relationships are statistically robust and biologically sensible, and multiple lines of experimental evidence support the predicted functional impact of both the well-known and the novel candidate drivers that are predicted by TCI. TCI provides a unified framework that integrates multiple types of SGAs and molecular phenotypes to estimate which genome perturbations are causally influencing one or more molecular/cellular phenotypes in an individual tumor. By identifying major candidate drivers and revealing their functional impact in an individual tumor, TCI sheds light on the disease mechanisms of that tumor, which can serve to advance our basic knowledge of cancer biology and to support precision oncology that provides tailored treatment of individual tumors.

# Identifying Tumor-specific Drivers

- Goal: identifying driver somatic genomic alterations (SGAs) of individual tumors
- A tumor usually hosts hundreds to thousands SGAs
  - Many passengers
  - Relatively few drivers
- Current knowledge of cancer driver genes is incomplete
  - 10 to 20% of tumors have no known drivers
  - 50% of tumors have  $\leq$  3 known drivers

# Tumor-specific Driver Identification (TDI) Algorithm\*



- Uses a bipartite graph representation
- Searches the graph for somatic genomic alterations (SGAs) that account for differentially expressed genes (DEGs) involved in an oncogenic process.
- Uses a Bayesian evaluation measure, which is tumor specific

*Hypothesis:* Those SGAs that predict DEGs well are good candidates as drivers of cancer

*Results:* Many known drivers recovered, several new candidate drivers discovered, many experimentally verified (6/6 perturbation confirmed).

# Common Scientific Purposes of Using of ACD on Empirical Data

#### 1. Finding alternatives .... to:

- 1. Theoretically specified causal hypotheses
- 2. Regression based causal inferences /

standard factor analysis latent variable models

#### 2. Generating candidate causes

### 3. Mechanism Discovery

# Mechanism Discovery

- 1. Educational Research
- 2. Brain Research
- 3. Climate
- 4. Miscellany (lots!)
- 5. Economics (for Kevin Hoover)

# **Educational Research**

College Plans - Wisconsin HS seniors, 1979, (N > 10,000)

Questions:

- 1. Does parental encouragement depend on sex, even controlling for SES and IQ
- 2. Does parental encouragement matter?

Measures:

- SES
- IQ
- SEX
- PE (parental encouragement
- CP (College Plans)





### **Online Course Data**

eb.cmu.edu

<	🐼 🏦 👗 https://odin.	.web.cmu.edu/jcou	rse/workbook/a 🚔	• 🕨 🖸 •	Google		Q			-	7
Calendar Philos	ophy 📄 Tetrad 📄 Causality	Lab 🛴 OLI - Cou	irse Entry 📩 ODIN	OLI Bla	ickboard	NAS-	VA »				
Suppose you ar your cell phone blace a call, an on END before f your call got	e traveling in a car . In the simulation d click on the "END you can try another through. Attempt a	and you w below, clic button to r call). The at least 10	ant to make k on the "SE end a call ( phone on th calls.	several c ND" butto you must e right wi	alls o n to click Il ring	n					
	SIMULATION OF CALL A	ATTEMPTS AND	CONNECTIONS			1	Ξ				
	ATTEMPTS: 1			1							
					RESET						
	> <b>Did I Get Th</b> Click Here	is?:		😻 t Eile	n <mark>ttps:/</mark> <u>E</u> dit	/odin.v View	web.cr Hi <u>s</u> tor	<b>nu.edu - Ce</b> y <u>B</u> ookmarks	ell Phor	n <b>e 1 - M</b> ; <u>H</u> elp	lozilla I
avascript:void openLgW	Did I Get Th Click Here	is?: esource.do?src=92	a573ffa80020c1500bd	287da26	n <mark>ttps:/</mark> Edit	<mark>/odin.v</mark> <u>V</u> iew	web.cr Hi <u>s</u> tor	<mark>mu.edu - Ce</mark> y <u>B</u> ookmarks	e <mark>ll Pho</mark> r s <u>T</u> ools	n <b>e 1 - M</b> ; <u>H</u> elp	ozilla I
iavascript:void openLgW	> <b>Did I Get Th</b> <u>Click Here</u> n(//jcourse/webui/resolver/link/r	is?: esource.do?src=9a	a573ffa80020c1500bd	287da26	nttps:/ <u>E</u> dit Cell	<mark>∕odin.v</mark> ⊻jew Phone	web.cr Histor	<mark>mu.edu - Ce</mark> y <u>B</u> ookmarks	e <mark>ll Pho</mark> i s <u>T</u> oole	ne 1 - M ; <u>H</u> elp	ozilla I
avascript:void openLgW	> <b>Did I Get Th</b> <u>Click Here</u> n('/jcourse/webui/resolver/link/r	is?: esource.do?src=90	a573ffa80020c1500bw	287da28 [] -	nttps:/ Edit Cell	/odin.v View Phone vestion	Histor Histor	nu.edu - Ce y <u>B</u> ookmarks estion 2 is question	s <u>T</u> ool	n <mark>e 1 - M</mark>	lozilla f
iavascript:void openLgW	> <b>Did I Get Th</b> <u>Click Here</u>	is?: esource.do?src=90	a573ffa80020r1500br	287da28	Cell	/odin.v vjew Phone estion cempt 1 Que Whick syste	Higtor Higtor 1 1   Qu for th stic	nu,edu - Ce y <u>B</u> ookmarks estion 2 is question on 2 he choices b	ell Phon s Iools	present	s the c
evescript:void openLgW	> <b>Did I Get Th</b> <u>Click Here</u>	is?: esource.do?src=9a	1573ffa80020c1500bd	2876426 [] -	Cell	Vodin.v Vjew Phone restion cempt 1 Que Whick syste	web.cr Higtor 1 1 1 Qu for th stic stic , {Ph	estion 2 is question con 2 ie choices b he simulation	est repon	present	s the c
Javascript:void openLgW	> <b>Did I Get Th</b> <u>Click Here</u>	is?: esource.do?src=90	a573ffa80020c1500bd	287da28 [ -	Edit Cell	Very View View Phane restion campt 1 Que Whidd syste O A O B	web.cr Higtor 1 for th stic for th stic for th for th for th stic for th stic	mu.edu - Ce y gookmarks estion 2 is question on 2 e choices b he simulation one Button one Button one Button	ell Phoi s Took	present ,end],	s the c
Jevascript:void openLgW	> <b>Did I Get Th</b> <u>Click Here</u>	is?: esource.do?src=9e	1573ffa90020c1500b0	287da26	titps:// Edit Cell	Vodin.v View Phone exempt 1 Que Whid syste O A O B	web.cr Higtor 1 1 for th stic for th mint , {Ph Atte {Ph Atte	nu.edu - Ce y Bookmarks estion 2 is question on 2 e choices b he simulation one Button one Button one Button i, no]}	ell Phores Tools	ne 1 - M s Help present ,end]}	ozilla (

- Causal Discovery Methods
  - → Learning Mechanism

### **15 Variables**

- Pre-test (%)
- Print-outs (% modules printed)
- Quiz Scores (avg. %)
- Voluntary Exercises (% completed)
- Final Exam (%)
- 9 other variables

#### Voluntary interaction → Learning: "The Doer Effect" 2005



Scheines, R., Leinhardt, G., Smith, J., and Cho, K. (2005) "Replacing Lecture with Web-Based Course Materials, *Journal of Educational Computing Research*, 32, 1, 1-26.



Koedinger, Kim, Jia, McLaughlin, & Bier (2015). Learning is not a spectator sport: Doing is better than watching for learning from a MOOC. In *Proceedings of the Second ACM Conference on Learning at Scale*.



#### "The Doer Effect" 2018

Figure 1. Using subsets of variables from the Georgia Tech Psychology dataset, three causal models were discovered using the PC algorithm and an alpha value of 0.05 as in [6].



Table 2. The causal model that was discovered on GTech's psychology dataset was estimated using data from datasets listed in the first row of the table.

	UMUC Biology	UMUC Info Sci	UMUC Psychology	UMUC Statistics	C@CM	UMUC Biology (sample)	UMUC Info Sci (sample)
#Students	3516	6112	89	61	383	300	300
Chi-square	78.89	18.44	1.04	28.33	14.30	11.92*	3.32*
DOF	2	2	2	2	4	2	2
P-value	0	0	0.59	0	0	0.02*	0.49*

Koedinger, K., Scheines, R., and Schaldebrand, P. (2018). "Is the Doer Effect Robust Across Multiple Data Sets?", *Proceedings of the 11th Intl. Conference on Educational Data Mining*, Buffalo, NY.

#### Is the Doer Effect Robust Across Multiple Data Sets?

Figure 2. Causal models of various datasets. To the bottom right of each model are the search algorithm and p-value cutoff for searching (alpha) used to discover the model. Below that are the model statistics when estimating the model on the dataset: Chi-square ( $\chi^2$ ), degrees of freedom in the model (DOF), and p-value.



Yes – in all samples doing was caually connected to success, and much more strongly than reading or watching so

Subsequent experiments confirm this finding – and indicate it is underestimated by 2x (unpublished)

Koedinger, K., Scheines, R., and Schaldebrand, P. (2018). "Is the Doer Effect Robust Across Multiple Data Sets?", *Proceedings of the 11th Intl. Conference on Educational Data Mining*, Buffalo, NY.

# Fractions Tutor (5<sup>th</sup> grade)



# Fluency vs. Understanding

• Understanding: sense-making processes



# Which to teach first?

• Fluency: fast, compiled, reliable



# Mediation Hypotheses



# Model Search Results

Understanding-1<sup>st</sup> reduces Fluency errors – which in turn increases post-test



# Model Search Results

Fluency-1<sup>st</sup> *increases* Understanding errors – which in turn *decreases* post-test



### **Educational Research**

Follow up Experiment: Sense-making first vs. Fluency-first

Sense-making first dramatically outperformed fluency-first condition on fraction learning

### Educational Data → Causal Discovery Methods → Instrumental Variable Detection → Causes of Educational Returns

#### Discovering Causal Models with Optimization: Confounders, Cycles, and Instrument Validity

Frederick Eberhardt Division of Humanities and Social Sciences, California Institute of Technology, fde@caltech.edu.

Nur Kaynar<sup>\*</sup> Samuel Curtis Johnson Graduate School of Management, Cornell University, sk2739@cornell.edu.

Auyon Siddiq Anderson School of Management, University of California, Los Angeles, auyon.siddiq@anderson.ucla.edu.

We propose a new optimization-based method for learning causal structures from observational data, a process known as *causal discovery*. Our method takes as input observational data over a set of variables and returns a graph in which causal relations are specified by directed edges. We consider a highly general search space that accommodates latent confounders and feedback cycles, which few extant methods do. We formulate the discovery problem as an integer program, and propose a solution technique that exploits the conditional independence structure in the data to identify promising edges for inclusion in the output graph. In the large-sample limit, our method recovers a graph that is (Markov) equivalent to the true datagenerating graph. Computationally, our method is competitive with the state-of-the-art, and can solve in minutes instances that are intractable for alternative causal discovery methods. We leverage our method to develop a graphical test for the validity of an instrumental variable and demonstrate it on the influential instruments for estimating the returns to education from Angrist and Krueger (1991) and Card (1993). In particular, our test complements existing instrument tests by revealing the precise causal pathways that undermine instrument validity, highlighting the unique merits of the graphical perspective on causality.

### Educational Data $\rightarrow$ Causal Discovery Methods $\rightarrow$ Instrumental Variable Detection $\rightarrow$ Causes of Educational Returns



Figure 6Edge frequencies over 50 bootstrap repetitions of EDGEGEN applied to Angrist and Krueger (1991) data for<br/>three levels of complexity penalty. Only edges with frequency  $\geq 0.1$  are shown.





#### Causal Discovery with Models: Behavior, Affect, and Learning in Cognitive Tutor Algebra

Stephen E. Fancsali Carnegie Learning, Inc.

Proceedings of the 7th International Conference on Educational Data Mining (EDM 2014)

#### ABSTRACT

Non-cognitive and behavioral phenomena, including gaming the system, off-task behavior, and affect, have proven to be important for understanding student learning outcomes. The nature of these phenomena requires investigations into their causal structure. For example, given that gaming the system has been associated with poorer learning outcomes, would reducing such behavior improve outcomes? Answering this question requires an understanding of whether gaming the system is a cause of poor outcomes, rather than, for example, only sharing a common cause with factors influencing learning. Because controlled experiments to settle such causal questions are often costly or impractical, we employ algorithmic search for the structure of graphical causal models from non-experimental data. Using sensor-free, data-driven detectors of behavior and affect, this work extends Baker and Yacef's notion of "discovery with models" to incorporate causal discovery and reasoning, resulting in an approach we call "causal discovery with models." We explore a case study of this approach using data from Carnegie Learning's Cognitive Tutor for Algebra and raise questions for future research.



# **Brain Research**

# fMRI $\rightarrow$ Causal Discovery Methods $\rightarrow$ Processing Mechanisms



#### Six problems for causal inference from fMRI

J.D. Ramsey <sup>a,\*</sup>, S.J. Hanson <sup>b</sup>, C. Hanson <sup>b</sup>, Y.O. Halchenko <sup>b</sup>, R.A. Poldrack <sup>c</sup>, C. Glymour <sup>d</sup>

\* Department of Philosophy, Carnegie Mellon University, Pittsburgh, PA 15213

b Department of Psychology, Rutgers University, Rumba Lab

<sup>c</sup> Imaging Research Center and Departments of Psychology and Neurobiology, University of Texas at Austin

<sup>d</sup> Department of Philosophy, Carnegie Mellon University, and Florida Institute for Human and Machine Cognition

#### ARTICLE INFO

Article history: Received 13 February 2009 Revised 7 August 2009 Accepted 31 August 2009 Available online 9 September 2009

#### ABSTRACT

Neuroimaging (e.g. fMRI) data are increasingly used to attempt to identify not only brain regions of interest (ROIs) that are especially active during perception, cognition, and action, but also the qualitative causal relations among activity in these regions (known as *effective connectivity*; Friston, 1994). Previous investigations and anatomical and physiological knowledge may somewhat constrain the possible hypotheses, but there often remains a vast space of possible causal structures. To find actual effective connectivity relations, search methods must accommodate indirect measurements of nonlinear time series dependencies, feedback, multiple subjects possibly varying in identified regions of interest, and unknown possible location-dependent variations in BOLD response delays. We describe combinations of procedures that under these conditions find feed-forward sub-structure characteristic of a group of subjects. The method is illustrated with an empirical data set and confirmed with simulations of time series of BOLD delays, with regions of interest missing at random for some subjects, measured with noise approximating the signal to noise ratio of the empirical data.

© 2009 Elsevier Inc. All rights reserved.



### fMRI (~44,000 voxels)





Causal network discovery



(ROI) ~10-20 Regions of Interest

# **Autism**

### Catherine Hanson, Rutgers

Subjects: Autistic Spectrum Disorder vs. Neurotypical

#### Usual Approach: Search for differential recruitment of brain regions

![](_page_59_Figure_4.jpeg)

### fMRI (~44,000 voxels)

![](_page_60_Picture_1.jpeg)

- Face processing task
- Theory of Mind task
- Action understanding task

![](_page_60_Picture_5.jpeg)

Causal network discovery

![](_page_60_Picture_7.jpeg)

(ROI) ~10-20 Regions of Interest

# Results

![](_page_61_Figure_1.jpeg)

# Causal Discovery $\rightarrow$ Autism

### Biwei Huang, CMU & UCSD

![](_page_62_Picture_2.jpeg)

Neural Engineering Techniques for Autism Spectrum Disorder Volume 1: Imaging and Signal Analysis

![](_page_62_Picture_4.jpeg)

2021, Pages 237-267

Chapter 12 - Diagnosis of autism spectrum disorder by causal influence strength learned from resting-state fMRI data

Biwei Huang

### Similar Follow on Work Catherine Hanson, Rutgers

![](_page_63_Picture_1.jpeg)

#### **HHS Public Access**

Author manuscript J Alcohol Drug Depend. Author manuscript; available in PMC 2017 October 11.

Published in final edited form as: J Alcohol Drug Depend. 2017 August; 5(4): . doi:10.4172/2329-6488.1000279.

Modeling Causal Relationships among Brain Areas in the Mesocorticolimbic System during Resting-State in Cocaine Users Utilizing a Graph Theoretic Approach

Suchismita Ray<sup>1,\*</sup>, Bharat B Biswal<sup>2</sup>, Ashley Aya<sup>1</sup>, Suril Gohel<sup>3</sup>, Aradhana Srinagesh<sup>1</sup>, Catherine Hanson<sup>4</sup>, and Stephen J. Hanson<sup>4</sup>

**Results**—The causal interaction pattern was different between the two groups. The feed-forward pattern found in cocaine smokers, between 7 ROIs of the MCLS during resting-state [ventral tegmental area (VTA)→hippocampus (HIPP)→ventral striatum (VenStri)→orbital frontal cortex (OFC), medial frontal cortex (MFC), anterior cingulate cortex (ACC), dorsolateral prefrontal cortex (DLPFC)], was absent in controls. That is, the subcortical VenStri area had a causal influence on four cortical brain areas only in cocaine users.

![](_page_63_Picture_8.jpeg)

Drug and Alcohol Dependence Volume 146, 1 January 2015, Page e77

![](_page_63_Picture_10.jpeg)

Modeling causal relationship between memory and craving-related brain networks in nontreatment seeking cocaine smokers using images, a graph theoretic approach

Suchismita Ray<sup>1</sup>, Catherine Hanson<sup>2</sup>, Margaret Haney<sup>3</sup>, Bharat Biswal<sup>4</sup>, Stephen J. Hanson<sup>25</sup>

# es-fMRI Data $\rightarrow$ Causal Discovery (FGES) $\rightarrow$ Emotion Network Structure (Amygdala)

bioRxiv preprint doi: https://doi.org/10.1101/214486; this version posted November 5, 2017. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

# Causal Mapping of Emotion Networks in the Human Brain: Framework and Preliminary Findings

Julien Dubois<sup>1,4</sup>, Hiroyuki Oya<sup>5</sup>, J. Michael Tyszka<sup>2</sup>, Matthew Howard III<sup>5</sup>, Frederick Eberhardt<sup>1</sup>, and Ralph Adolphs<sup>1,2,3</sup>

<sup>1</sup>Division of Humanities and Social Sciences,
 <sup>2</sup>Division of Biology, and
 <sup>3</sup>Chen Neuroscience Institute, California Institute of Technology, Pasadena CA 91125, USA
 <sup>4</sup>Department of Neurosurgery, Cedars-Sinai Medical Center, Los Angeles, CA 90048, USA
 <sup>5</sup>Department of Neurosurgery, Human Brain Research Laboratory, University of Iowa, IA 52241, USA

# es-fMRI Data $\rightarrow$ Causal Discovery (FGES) $\rightarrow$ Emotion Network Structure (Amygdala)

![](_page_65_Figure_1.jpeg)

# **Climate Research**

### **Climate Teleconnections**

Journal of Machine Learning Research 9 (2008) 967-991

Submitted 9/06; Revised 9/07; Published 5/08

#### Search for Additive Nonlinear Time Series Causal Models

Tianjiao Chu

TIC19@PITT.EDU

Department of Obstetrics, Gynecology & Reproductive Sciences University of Pittsburgh 204 Craft Ave., Room B409 Pittsburgh, PA 15213, USA

CG09@ANDREW.CMU.EDU

Clark Glymour Department of Philosophy Carnegie Mellon University Pittsburgh, PA 15213, USA

Editor: Greg Ridgeway

#### Abstract

Pointwise consistent, feasible procedures for estimating contemporaneous linear causal structure from time series data have been developed using multiple conditional independence tests, but no such procedures are available for non-linear systems. We describe a feasible procedure for learning a class of non-linear time series structures, which we call additive non-linear time series. We show that for data generated from stationary models of this type, two classes of conditional independence relations among time series variables and their lags can be tested efficiently and consistently using tests based on additive model regression. Combining results of statistical tests for these two classes of conditional independence relations and the temporal structure of time series data, a new consistent model specification procedure is able to extract relatively detailed causal information. We investigate the finite sample behavior of the procedure through simulation, and illustrate the application of this method through analysis of the possible causal connections among four ocean indices. Several variants of the procedure are also discussed.

Keywords: conditional independence test, contemporaneous causation, additive model regression, Granger causality, ocean indices

### **Climate Teleconnections**

Our data set consists of the following 4 ocean climate indices, recorded monthly from 1958 to 1999, each forming a time series of 504 time steps:

- SOI Southern Oscillation Index: Sea Level Pressure (SLP) anomalies between Darwin and Tahiti
- WP Western Pacific: Low frequency temporal function of the 'zonal dipole' SLP spatial pattern over the North Pacific.
- AO Arctic Oscillation: First principal component of SLP poleward of 20° N
- NAO North Atlantic Oscillation: Normalized SLP differences between Ponta Delgada, Azores and Stykkisholmur, Iceland

![](_page_67_Figure_6.jpeg)

![](_page_67_Figure_7.jpeg)

# **Climate Research**

### Time Series Climate Data $\rightarrow$ Causal Discovery Methods $\rightarrow$ Causal Mechanisms in Earth's Climate

#### nature communications

Explore content V About the journal V Publish with us V

nature > nature communications > perspectives > article

Perspective Open Access Published: 14 June 2019

# Inferring causation from time series in Earth system sciences

Jakob Runge <sup>[2]</sup>, <u>Sebastian Bathiany</u>, <u>Erik Bollt</u>, <u>Gustau Camps-Valls</u>, <u>Dim Coumou</u>, <u>Ethan Deyle</u>, <u>Clark</u> <u>Glymour</u>, <u>Marlene Kretschmer</u>, <u>Miguel D. Mahecha</u>, <u>Jordi Muñoz-Marí</u>, <u>Egbert H. van Nes</u>, <u>Jonas Peters</u>, <u>Rick Quax</u>, <u>Markus Reichstein</u>, <u>Marten Scheffer</u>, <u>Bernhard Schölkopf</u>, <u>Peter Spirtes</u>, <u>George Sugihara</u>, <u>Jie</u> <u>Sun</u>, <u>Kun Zhang</u> & <u>Jakob Zscheischler</u>

Nature Communications 10, Article number: 2553 (2019) Cite this article

### Time Series Climate Data $\rightarrow$ Causal Discovery Methods → Causal Structure of the Climate

![](_page_69_Picture_1.jpeg)

	Jakob Runge						
	German Aerospace Cent Verified email at dlr.de - J Causal Inference Time	<u>er</u> and T <u>Homepa</u> Series	U Berlin g <u>e</u> Statistical Machine Learning	Information Theory	Earth Sc	ience	
TITLE				CITE	D BY	YEAR	
Inferring causation J Runge, S Bathiany, E Nature communications	from time series in Earth Bollt, G Camps-Valls, D Coum 10 (1), 2553	system iou, E De	sciences yle,		501	2019	
Detecting and quan J Runge, P Nowack, M Science advances 5 (11	<b>tifying causal associatio</b> Kretschmer, S Flaxman, D Sej ), eaau4996	ns in lar idinovic	ge nonlinear time series data	isets	481	2019	
Escaping the curse J Runge, J Heitzig, V Pe Physical review letters 1	of dimensionality in estination toukhov, J Kurths 08 (25), 258701	mating r	nultivariate transfer entropy		314	2012	
Causal network rec estimation J Runge Chaos: An Interdisciplin	onstruction from time se	ries: Fro	om theoretical assumptions to	o practical	269	2018	
Identifying causal g J Runge, V Petoukhov, Nature communications	ateways and mediators i JF Donges, J Hlinka, N Jajcay 6 (1), 8502	n comp , M Vejme	lex spatio-temporal systems elka,		246	2015	
Using causal effect M Kretschmer, D Coum Journal of climate 29 (1)	networks to analyze difference ou, JF Donges, J Runge I), 4069-4081	erent Ar	ctic drivers of midlatitude win	ter circulation	235	2016	
Quantifying the stree correlation and a no J Runge, V Petoukhov, Journal of climate 27 (2)	ngth and delay of climat wel measure based on g J Kurths , 720-739	ic intera graphica	ctions: The ambiguities of cro Il models	DSS	172	2014	

### Time Series Climate Data $\rightarrow$ Causal Discovery Methods $\rightarrow$ Causal Structure of the Climate

![](_page_70_Picture_1.jpeg)

TITLE

Jakob Rung	×	FOLLOW			
<u>German Aerospac</u> Verified email at d	<u>e Center</u> and T Ir.de - <u>Homepa</u>	ΓU Berlin g <u>e</u>			
Causal Inference	Time Series	Statistical Machine Learning	Information Theory	Earth Sc	ience
			CITE	D BY	YEAR

Inferring causation from time series in Earth system sciences J Runge, S Bathiany, E Bollt, G Camps-Valls, D Coumou, E Deyle, Nature communications 10 (1), 2553	501	2019
Detecting and quantifying causal associations in large nonlinear time series datasets J Runge, P Nowack, M Kretschmer, S Flaxman, D Sejdinovic Science advances 5 (11) eaau4996	481	2019

![](_page_70_Figure_4.jpeg)

Fig. 2. Source: Fig.2 of Runge et al. (2019) on the higher detection power of PCMCI partial correlation compared to Correlation, and FullCI partial correlation.

![](_page_71_Picture_0.jpeg)

Energy Reports Volume 7, November 2021, Pages 6196-6204

![](_page_71_Picture_2.jpeg)

#### Research paper

#### Exploring nonlinearity on the CO<sub>2</sub> emissions, economic production and energy use nexus: A causal discovery approach

#### Peter Martey Addo <sup>a</sup> o , <u>Christelle Manibialoa</u> <sup>a</sup>, <u>Florent McIsaac</u> <sup>b</sup>

#### Show more 🗸

🕂 Add to Mendeley 😪 Share 🍠 Cite

https://doi.org/10.1016/j.egyr.2021.09.026 🏼 🛪	Get rights and content 🛛
Under a Creative Commons license	<ul> <li>open access</li> </ul>

#### Highlights

- At the global scale, energy is central in the economic and decarbonization debate.
- Coordination between countries will improve policy effectiveness on climate actions.
- Energy and climate policies in three regions will impact global economic dynamics.
- Examining interactions between carbon emissions, economic production, and energy use.
- Too rapid a transition to net-zero emissions in the energy sector may have negative consequences for economic growth.

![](_page_71_Figure_15.jpeg)

#### Download : Download high-res image (93KB) Download : Download full-size image

Fig. 5. Causal dependence at global level. The PCMCI with CMIknn independence test produces better predictive power with NRMSE = 0.77 compared to Parcorr and GPDC.

![](_page_71_Figure_18.jpeg)

Fig. 7. Causal dependence of all regions. The PCMCI with GPDC indeper


A Causality-Based View of the Interaction between Synoptic- and Planetary-Scale Atmospheric Disturbances

Savini M. Samarasinghe, Yi Deng, and Imme Ebert-Uphoff

Online Publication: 01 Feb 2020

Print Publication: 01 Mar 2020

DOI: https://doi.org/10.1175/JAS-D-18-0163.1

Page(s): 925-941

Article History Download PDF © Get Permissions

#### Abstract/Excerpt Full Text PDF

#### Abstract

This paper reports preliminary yet encouraging findings on the use of causal discovery methods to understand the interaction between atmospheric planetary- and synopticscale disturbances in the Northern Hemisphere. Specifically, constraint-based structure learning of probabilistic graphical models is applied to the spherical harmonics decomposition of the daily 500-hPa geopotential height field in boreal winter for the period 1948–2015. Active causal pathways among different spherical harmonics components are identified and documented in the form of a temporal probabilistic graphical model. Since, by definition, the structure learning algorithm used here only robustly identifies linear causal effects, we report only causal pathways between two groups of disturbances with sufficiently large differences in temporal and/or spatial scales, that is, planetary-scale (mainly zonal wavenumbers 1–3) and synoptic-scale disturbances (mainly zonal wavenumbers 6–8). Daily reconstruction of geopotential heights using only interacting scales suggest that the modulation of synoptic-scale disturbances by planetary-scale disturbances is best characterized by the flow of information from a zonal wavenumber-1 disturbance to a synoptic-scale circumglobal wave train whose amplitude peaks at the North Pacific and North Atlantic storm-track region. The feedback of synoptic-scale to planetary-scale disturbances manifests itself as a zonal wavenumber-2 structure driven by synoptic-eddy momentum fluxes. This wavenumber-2 structure locally enhances the East Asian trough and western Europe ridge of the wavenumber-1 planetary-scale disturbance that actively modulates the activity of synoptic-scale disturbances. The winter-mean amplitude of the actively interacting disturbances are characterized by pronounced fluctuations across interannual to decadal time scales.

# Biology

Pacific Symposium on Biocomputing 7:498-509 (2002)

### DISCOVERY OF CAUSAL RELATIONSHIPS IN A GENE-REGULATION PATHWAY FROM A MIXTURE OF EXPERIMENTAL AND OBSERVATIONAL DNA MICROARRAY DATA

C. YOO\*, V. THORSSON<sup>\*</sup>, and G.F. COOPER\* \*Center for Biomedical Informatics, University of Pittsburgh 8084 Forbes Tower, 200 Lothrop St., Pittsburgh PA 15213 \*The Institute for Systems Biology

4225 Roosevelt Way NE, Suite 200, Seattle Washington 98105

This paper reports the methods and results of a computer-based search for causal relationships in the gene-regulation pathway of galactose metabolism in the yeast *Saccharomyces cerevisiae*. The search uses recently published data from cDNA microarray experiments. A Bayesian method was applied to learn causal networks from a mixture of observational and experimental gene-expression data. The observational data were gene-expression levels obtained from unmanipulated "wild-type" cells. The experimental data were produced by deleting ("knocking out") genes and observing the expression levels of other genes. Causal relations predicted from the analysis on 36 galactose gene pairs are reported and compared with the known galactose pathway. Additional exploratory analyses are also reported.

### Discovering Signaling Pathways – Karen Sachs



Figure 10. Overview of Markov neighborhood algorithm. A set of preliminary experiments are used to determine variable neighborhoods; these include variables which appear to be reasonable candidate parents based on a dependence metric (e.g. correlation). Sets of m molecules from each neighborhood are profiled under various stimulus conditions, and a constrained Bayesian network structure learning analysis is performed, in which candidate variable parents are selected from the (profiled portion of the) variable's neighborhood.

### Discovering Signaling Pathways – Karen Sachs



Figure 7. Model results contrasted to known influence connections. Literature reports were used to assess the retrieved network. "Classic/Expected" edges indicate edges that are well established in the literature, "Reported" edges are not well established and have never been reported in T-cells, "Reversed" indicates an edge which is directed incorrectly, and "Missed" indicates well established edges that were not found by the modeling approach.

### **Biology and Health**

### Psychological Medicine

cambridge.org/psm

### **Original Article**

\*Drs. Kummerfeld and Vinogradov contributed equally as senior authors to this manuscript.

Cite this article: Miley K, Meyer-Kalos P, Ma S, Bond DJ, Kummerfeld E, Vinogradov S (2023). Causal pathways to social and occupational functioning in the first episode of schizophrenia: uncovering unmet treatment needs. *Psychological Medicine* **53**, 2041–2049. https://doi.org/10.1017/S0033291721003780

Received: 10 April 2021 Revised: 30 August 2021 Accepted: 1 September 2021 First published online: 8 October 2021

#### Key words:

Early schizophrenia; motivation; social cognition; causal discovery; functional outcomes; machine learning

#### Author for correspondence:

Kathleen Miley, E-mail: mile0087@umn.edu

Causal pathways to social and occupational functioning in the first episode of schizophrenia: uncovering unmet treatment needs

Kathleen Miley<sup>1,2</sup>, Piper Meyer-Kalos<sup>1</sup>, Sisi Ma<sup>3</sup>, David J. Bond<sup>1</sup>, Erich Kummerfeld<sup>3,\*</sup> and Sophia Vinogradov<sup>1,\*</sup>

<sup>1</sup>Department of Psychiatry and Behavioral Sciences, University of Minnesota, Minneapolis, MN, USA; <sup>2</sup>School of Nursing, University of Minnesota, Minneapolis, MN, USA and <sup>3</sup>Institute for Health Informatics, University of Minnesota, Minneapolis, MN, USA

#### Abstract

**Background.** We aimed to identify unmet treatment needs for improving social and occupational functioning in early schizophrenia using a data-driven causal discovery analysis. **Methods.** Demographic, clinical, and psychosocial measures were obtained for 276 participants from the Recovery After an Initial Schizophrenia Episode Early Treatment Program (RAISE-ETP) trial at baseline and 6-months, along with measures of social and occupational functioning from the Quality of Life Scale. The Greedy Fast Causal Inference algorithm was used to learn a partial ancestral graph modeling causal relationships across baseline variables and 6-month functioning. Effect sizes were estimated using a structural equation model. Results were validated in an independent dataset (N = 187).

**Results.** In the data-generated model, greater baseline socio-affective capacity was a cause of greater baseline motivation [Effect size (ES) = 0.77], and motivation was a cause of greater baseline social and occupational functioning (ES = 1.5 and 0.96, respectively), which in turn were causes of their own 6-month outcomes. Six-month motivation was also identified as a cause of occupational functioning (ES = 0.92). Cognitive impairment and duration of untreated psychosis were not direct causes of functioning at either timepoint. The graph for the validation dataset was less determinate, but otherwise supported the findings.

**Conclusions.** In our data-generated model, baseline socio-affective capacity and motivation are the most direct causes of occupational and social functioning 6 months after entering treatment in early schizophrenia. These findings indicate that socio-affective abilities and motivation are specific high-impact treatment needs that must be addressed in order to promote optimal social and occupational recovery.

Psychological Medicine

cambridge.org/psm

Causal pathways to social and occupational functioning in the first episode of schizophrenia: uncovering unmet treatment needs

#### **Original Article**

\*Drs. Kummerfeld and Vinogradov contributed equally as senior authors to this manuscript.

Kathleen Miley<sup>1,2</sup> , Piper Meyer-Kalos<sup>1</sup>, Sisi Ma<sup>3</sup>, David J. Bond<sup>1</sup>, Erich Kummerfeld<sup>3,\*</sup> and Sophia Vinogradov<sup>1,\*</sup>



Fig. 2. Study 1 complete PAG and functional outcome subgraph. All variables represent baseline unless specified as 6M (6 month). Alcohol use = days of alcohol use



Fig. 4. Study 2: validation complete PAG and functional outcome subgraph. All variables represent baseline unless specified as 6M (6-month). Antipsychotic years = unarc of antipoychotic user Collected Clebal Improvements Scale for Schippehrenia, Eacial Affect Recognition =



Volume 28, Issue 1 January 2012

#### **Article Contents**

Abstract

1 INTRODUCTION

2 METHODS

3 RESULTS

4 DISCUSSION

5 CONCLUSION

ACKNOWLEDGEMENTS

REFERENCES

Author notes

Supplementary data

< Previous Next >

#### JOURNAL ARTICLE

# Inferring gene regulatory networks from gene expression data by path consistency algorithm based on conditional mutual information **a**

Xiujun Zhang, Xing-Ming Zhao, Kun He, Le Lu, Yongwei Cao, Jingdong Liu, Jin-Kao Hao, Zhi-Ping Liu ⊠, Luonan Chen ⊠ Author Notes

Bioinformatics, Volume 28, Issue 1, January 2012, Pages 98–104, https://doi.org/10.1093/bioinformatics/btr626 Published: 15 November 2011 Article history ▼

🕨 PDF 📲 Split View 😘 Cite 🔑 Permissions < Share 🔻

#### Abstract

Motivation: Reconstruction of gene regulatory networks (GRNs), which explicitly represent the causality of developmental or regulatory process, is of utmost interest and has become a challenging computational problem for understanding the complex regulatory mechanisms in cellular systems. However, all existing methods of inferring GRNs from gene expression profiles have their strengths and weaknesses. In particular, many properties of GRNs, such as topology sparseness and non-linear dependence, are generally in regulation mechanism but seldom are taken into account simultaneously in one computational method.

**Results:** In this work, we present a novel method for inferring GRNs from gene expression data considering the non-linear dependence and topological structure of GRNs by employing path consistency algorithm (PCA) based on conditional mutual information (CMI). In this algorithm, the conditional dependence between a pair of genes is represented by the CMI between them. With the general hypothesis of Gaussian distribution underlying gene expression data, CMI between a pair of genes is computed by a concise formula involving the covariance matrices of the related gene expression profiles. The method is validated on the benchmark GRNs from the DREAM challenge and the widely used SOS DNA repair network in *Escherichia coli*. The cross-validation results confirmed the effectiveness of our method. Besides its high accuracy, our method is able to distinguish direct (or causal) interactions from indirect associations.

### SCIENTIFIC REPORTS

natureresearch

# OPEN Challenges and Opportunities with Causal Discovery Algorithms: Application to Alzheimer's Pathophysiology

Xinpeng Shen<sup>1\*</sup>, Sisi Ma<sup>1</sup>, Prashanthi Vemuri<sup>2</sup>, Gyorgy Simon<sup>1\*</sup> & the Alzheimer's Disease Neuroimaging Initiative<sup>†</sup>

Causal Structure Discovery (CSD) is the problem of identifying causal relationships from large quantities of data through computational methods. With the limited ability of traditional association-based computational methods to discover causal relationships, CSD methodologies are gaining popularity. The goal of the study was to systematically examine whether (i) CSD methods can discover the known causal relationships from observational clinical data and (ii) to offer guidance to accurately discover known causal relationships. We used Alzheimer's disease (AD), a complex progressive disease, as a model because the well-established evidence provides a "gold-standard" causal graph for evaluation. We evaluated two CSD methods, Fast Causal Inference (FCI) and Fast Greedy Equivalence Search (FGES) in their ability to discover this structure from data collected by the Alzheimer's Disease Neuroimaging Initiative (ADNI). We used structural equation models (which is not designed for CSD) as control. We applied these methods under three scenarios defined by increasing amounts of background knowledge provided to the methods. The methods were evaluated by comparing the resulting causal relationships with the "gold standard" graph that was constructed from literature. Dedicated CSD methods managed to discover graphs that nearly coincided with the gold standard. For best results, CSD algorithms should be used with longitudinal data providing as much prior knowledge as possible.

## Summary / Conclusions

A lot of good work on important scientific questions -

- But distributed across disciplines, so largely unknown
- Almost all of it by algorithm builders/developers
- Too little of it includes experimental follow-up

By far the most common use is mechanism discovery

We need a sub-discipline devoted to applications