Causal Discovery by Computer

Clark Glymour

Carnegie Mellon University

Outline

- 1. A century of mistakes about causation and discovery:
 - 1. Fisher
 - 2. Yule
 - 3. Spearman/Thurstone
- 2. Search for causes is statistical estimation
- 3. Strategies and algorithms
- 4. Examples
- 5. Prospects: The Center for Causal Discovery

A Conversation with R.A. Fisher

Science is about causation. Causes can only be discovered by experiment.





Physics and astronomy. Then statistics and biology, especially evolution.



Say again!?

Continued...

Very well: Science is about causation. Causes can only be discovered by experiment.







So much then for my discovery of the cause of the motions of the planets, and the cause of the tides.

> So much as well for my discovery of the origin of species.



Lung Cancer is associated with smoking. The association can equally be explained in three ways:

- 1. Early undetected cancers cause smoking
- 2. Genetics causes both smoking and cancer
- *3. Smoking causes cancer The data cannot distinguish these explanations.*

The argument is quite general: Associations cannot distinguish cause from effect from joint effects of the unobserved

Even Genius Makes Mistakes



Suppose Data:

- Parents Smoking and Smoking are not associated with Job
- Parents Smoking is not associated with Cancer
- Smoking and Cancer are associated
- Parents Smoking and Smoking are associated
- Job and Cancer are associated

Suppose Data:

- Parents Smoking and Smoking are not associated with Job
- Parents Smoking is not associated with Cancer
- Smoking and Cancer are associated
- Parents Smoking and Smoking are associated
- Job and Cancer are associated

Then every causal explanation of the data entails that



Another example:

Suppose data were to show:

- z X and Y are not associated
- z Z is associated with X and Y



- z X and Y *are* associated conditional on a value of Z
- z X, Y, Z are associated with R
- z X, Y are not associated with R conditional on Z

Then Z is a cause of R and there are no confounding common causes of Z and R.

George Udney Yule and Regression



Yule championed multiple regression as a means to estimate causal effects. He was aware of some, but not all, of its foibles: Unmeasured Unmeasured $\beta_{xz,y}$ 7

Although X and Z have no unmeasured common cause, and X has no influence on Z, regression of Z on X,Y results in a

9

Spearman, Thurstone and Factor Analysis



Thurstone's Equivocation: "Vectors of Mind" = Linear reduction of correlations





L.L. Ohurston

Upshot for the 20th Century

- z Causal inferences should not be made from nonexperimental data
- z But if you must, use regression
- z Or factor analysis
- z Or "potential outcomes"

Fear of Search





Terry Speed: There may be better search methods than regression, but only regression should be used.

~ 1990: Rethinking:

- z Relationships between causal hypotheses and probability hypotheses
- z Search

Search for causal relations is statistical estimation

Hynothesis Snace		A	В	
Hypothesis opacei	А	?	?	?
	В	?	?	?
	С	?	?	?

Data: A II C, A
$$\underline{H}$$
 B, B \underline{H} C, A \underline{H} C | B

Statistical inference:

A -> B <- C:

	Α	В	
А	0	1	0
В	0	0	1
С	0	1	0

What's the Difference?

In conventional estimation we are estimating a *probability distribution*—an unobserved distribution covering present and potential future observations.

In causal estimation, we are estimating both a current probability distribution *and* the *probability distributions that would result from various interventions.*

Interventions Change Probabilities



Observed

Force a value on Y; distribution of X doesn't change. Dependence of Y on X is broken

Force a value on X; distribution of Y changes

Goals are the same as in Experimental Design

<u>First Goal:</u> To estimate whether there are **some** values of other variables **C** such that **some** intervention that changes the distribution of A will change the distribution of B when **C** variables are forced to have those values.

<u>Secondary Goals</u>: to estimate the signs or strengths of effects.



Conventional Statistical Estimation Has Assumptions, E.G.

- z Distribution Family?
- z I.I.D Sampling?
- z Censored Data?
- z Stationary Time Series?
- z Variance Known?
- z Prior Probabilities?
- z Likelihood Function (i.e., the "model")

Estimation Are Standards for Search

- z Asymptotic convergence to true information under explicit assumptions.
- z Unbiased (Expected value is the true value)
- z Finite sample error probabilities.
- z Tests of assumptions.
- z Robustness to "small" violations of assumptions.

Representation

Z Represent causal relations on a collection
V of variables by directed edges X -> Y

Z A directed edge indicates that some intervention on X would change the distribution of Y if all other variables in V (that are not effects of X) were held constant at some values.

What Extra Assumptions Are Needed for Causal Inference?

Answer: Generalizations of principles of experimental design:

- Markov Condition
- Faithfulness Condition

Experimental Principles

Why do we randomize in experiments?

Because we expect (or hope) that on average randomizing X will remove any association between X and Y due to common causes of X and Y.

Randomize X Y And *thus*, the association (or it's absence) of X and Y in the experiment will measure the effect (or it's

absence) of X on Y.

Why in experiments do we "stratify" values of potential confounders, Z, (i.e. arrange subsamples in which Z is constant)? One reason: Because we expect that the association of X and Y when Z is constant (i.e., conditioned on) will measure the effect of X on Y.

Experimental Principles: Blinding

When we think Y is an effect of X and we do not know whether Z is an effect of Y, and we want to know whether X causes Z



- z We do **not** force Y to be constant when we randomize X.
- z We do **not** condition on Y when we randomize X.
- z But if we want to know if there is a direct influence of X on Z, we **do** condition on Y.

Generalizations

- Z Markov: X is independent of all other variables (except for effects of X) conditional on the direct causes of X.
- Z Faithfulness: All independencies and conditional independencies among variables in a system of variables follow from the Markov condition for the true causal graph of the system.

A Flow of Mathematical Results Since 1990

- **z** *Sufficient* conditions for recovering causal information:
- z Markov
- z Acyclic
- z Faithfulness
- z I.I.D sampling
- z No unrecorded common causes

And a host of alternative sufficient conditions.

Examples for which There Are Correct Search Procedures

U1 \longrightarrow U2 \leftarrow U3 \swarrow \checkmark \checkmark \checkmark \checkmark \checkmark \checkmark \checkmark Linearity of X variables X 1 X2 X3 X4 X5 X6

X1 \rightarrow X2 X1 \rightarrow X2 \rightarrow X3 X1 \rightarrow X2 \rightarrow X3 \downarrow Linear X1 \leftrightarrow X2 \downarrow X3 \leftrightarrow X4 \downarrow X1 X2 X3 X4 \downarrow U1 \rightarrow U2 \downarrow U2

The Zen of Search: Don't Look (where you don't need to)



Applications

In flight recalibration Acid Rain





College Dropouts



Lead & IQ

Gene Regulators

Autism Spectrum







Center for Causal Discovery

- 1. New Algorithms
- 2. Moving from small discovery problems to huge problems
- 3. Making software easily usable
- 4. Training
- 5. Science Applications:
 - V Cellular pathways in breast cancer
 - yin lung disease
 - Brain processing in normals and autistics

Cancer Signaling and Big Data



z Genomic data

- y Somatic mutations
- y Somatic copy number alterations
- y DNA methylation
- z Proteomic data
 - y Quantity of certain signaling proteins
 - y Chemical modification (phosphorylation) of proteins
- z Transcriptomic data
 - y mRNA
 - y miRNA
 - y IncRNA

Signaling Pathway and Causal Network



Identify associations between clinical features, omics and lung disease using graphical models (undirected)





© 2015 Benos lab / Univ of Pittsburgh

Finding the Neural Mechanisms in Autism









Deep Boltzmann Machine



© 2015 Benos lab / Univ of Pittsburgh

Comparison of MGM-Learn flavors to other structure learning methods



© 2015 Benos lab / Univ of Pittsburgh

Inferring Unobserved Intermediate Causes—Between Somatic Gene Anomalies and mRNA Transcripts

Protein Signalling Network



Hippocampal Pathways



Identification of Communicating Sub-Regions of the Hippocampus









Distinguishing Autistics by Graphical Search with fMRI



Predicted	Actual	
	NT	AT
NT	.92	.01
AT	.08	.99

The Big Task for Big Data

Z Showing that we can use big biomedical data to discover novel, important causal relations that can be experimentally confirmed.