Simplicity of Additive Noise Models

Jonas Peters ETH Zürich - Marie Curie (IEF)

Workshop on Simplicity and Causal Discovery Carnegie Mellon University

7th June 2014





• ETH Zürich:

Peter Bühlmann, Jan Ernest

- Max-Planck-Institute Tübingen: Dominik Janzing, Bernhard Schölkopf
- University of Amsterdam, Radboud University Nijmegen: Joris Mooij
- UC Berkeley:

Sivaraman Balakrishnan, Martin Wainwright



Assume $P(X_1, \ldots, X_4)$ has been generated by



Can the directed acyclic graph be recovered from $P(X_1, \ldots, X_4)$?

Assume $P(X_1, \ldots, X_4)$ has been generated by



Can the directed acyclic graph be recovered from $P(X_1, \ldots, X_4)$? No.

Proposition

Given a distribution P, we can find a SEM generating P for each graph G, such that P is Markov with respect to G.

JP: Restricted Structural Equation Models for Causal Inference, PhD Thesis 2012 (and probably others?)

Proposition

Given a distribution P, we can find a SEM generating P for each graph G, such that P is Markov with respect to G a lot of graphs.

JP: Restricted Structural Equation Models for Causal Inference, PhD Thesis 2012 (and probably others?)

Assume $P(X_1, \ldots, X_4)$ has been generated by



Additive Noise Models.

Assume $P(X_1, \ldots, X_4)$ has been generated by



Additive Noise Models with Gaussian noise.

Assume $P(X_1, \ldots, X_4)$ has been generated by



Additive Noise Models with Gaussian noise. Can the DAG be recovered from $P(X_1, ..., X_4)$?

Assume $P(X_1, \ldots, X_4)$ has been generated by



Additive Noise Models with Gaussian noise. Can the DAG be recovered from $P(X_1, ..., X_4)$? Yes iff f_i nonlinear. JP, J. Mooij, D. Janzing and B. Schölkopf: *Causal Discovery with Continuous Additive Noise Models*, to appear in JMLR

Jonas Peters (ETH Zurich)

Simplicity of Additive Noise Models









GAUL GAUSS "the LINEAR"

Proposition

Given a Gaussian distribution P, we can find a linear Gaussian SEM generating P for each graph G, such that P is Markov with respect to G.

A. Hauser: Causal Inference from Interventional Data, PhD Thesis 2013 (and probably others?)



Consider a distribution generated by

$$Y = f(X) + N_Y$$

with $N_Y, X \stackrel{ind}{\sim} \mathcal{N}$



Consider a distribution generated by

$$Y = f(X) + N_Y$$

with $N_Y, X \stackrel{ind}{\sim} \mathcal{N}$



Then, if f is nonlinear, there is no



JP, J. Mooij, D. Janzing and B. Schölkopf: Causal Discovery with Continuous Additive Noise Models, to appear in JMLR

Consider a distribution corresponding to

$$Y = X^3 + N_Y$$

with $N_Y, X \stackrel{ind}{\sim} \mathcal{N}$

$$(X) \longrightarrow (Y)$$

with

 $X \sim \mathcal{N}(1, 0.5^2)$ $N_Y \sim \mathcal{N}(0, 0.4^2)$











Method No. 1: Testing for independent residuals

Regress each variable on the other and check whether the residuals are independent of the input/independent/explanatory variable.





Jonas Peters (ETH Zurich)

Simplicity of Additive Noise Models

7th June 2014



F. H. Messerli: Chocolate Consumption, Cognitive Function, and Nobel Laureates, N Engl J Med 2012



F. H. Messerli: Chocolate Consumption, Cognitive Function, and Nobel Laureates, N Engl J Med 2012





No (not enough) data for chocolate



No (not enough) data for chocolate



... but we have data for coffee!



Correlation: 0.698 *p*-value: $< 2.2 \cdot 10^{-16}$



Correlation: 0.698 *p*-value: $< 2.2 \cdot 10^{-16}$

Coffee \rightarrow Nobel Prize: Dependent residuals (*p*-value of $5.1 \cdot 10^{-78}$). Nobel Prize \rightarrow Coffee: Dependent residuals (*p*-value of $3.1 \cdot 10^{-12}$).

 \Rightarrow Model class too small? Causally insufficient?



Correlation: 0.698 *p*-value: $< 2.2 \cdot 10^{-16}$

Coffee \rightarrow Nobel Prize: Dependent residuals (*p*-value of $5.1 \cdot 10^{-78}$). Nobel Prize \rightarrow Coffee: Dependent residuals (*p*-value of $3.1 \cdot 10^{-12}$).

 \Rightarrow Model class too small? Causally insufficient? Question: When is a *p*-value too small?

Method No. 1: Testing for independent residuals

Regress each variable on the other and check whether the residuals are independent of the input/independent/explanatory variable.



Method No. 1: Testing for independent residuals

Regress each variable on the other and check whether the residuals are independent of the input/independent/explanatory variable.

୫ ଅନ୍

 \rightsquigarrow translates to p > 2 variables.

Nice: correct (in population), no distributional assumption about noise. Problem: does not scale well to large p although it's $\mathcal{O}(p^2)$ ind. tests.



 \succ






Method No. 2: Minimizing KL

Estimate the direction that corresponds to the closest subspace (details follow).



Proposition

Assume P(X, Y) is generated by

$$Y = \beta X^2 + N_Y$$

with independent $X \sim \mathcal{N}(0, \sigma_X^2)$ and $N_Y \sim \mathcal{N}(0, \sigma_{N_Y}^2)$.

Proposition

Assume P(X, Y) is generated by

$$Y = \beta X^2 + N_Y$$

with independent $X \sim \mathcal{N}(0, \sigma_X^2)$ and $N_Y \sim \mathcal{N}(0, \sigma_{N_Y}^2)$.

Then

$$\inf_{Q\in\{Q:Y\to X\}} \operatorname{KL}(P \mid\mid Q) > 0 \qquad \text{if } \beta \neq 0 \,.$$

Proposition

Assume P(X, Y) is generated by

$$Y = \beta X^2 + N_Y$$

with independent $X \sim \mathcal{N}(0, \sigma_X^2)$ and $N_Y \sim \mathcal{N}(0, \sigma_{N_Y}^2)$.

Then

$$\inf_{Q \in \{Q: Y \to X\}} \operatorname{KL}(P \mid\mid Q) = \frac{1}{2} \log \left(1 + 2\beta^2 \frac{\sigma_X^4}{\sigma_{N_Y}^2} \right)$$

Proposition

Assume P(X, Y) is generated by

 $Y = f(X) + N_Y$

with independent $X \sim \mathcal{N}(0, \sigma_X^2)$ and $N_Y \sim \mathcal{N}(0, \sigma_{N_Y}^2)$.

Then

$$\inf_{Q \in \{Q: Y \to X\}} \operatorname{KL}(P \mid\mid Q) \geq$$

Proposition

Assume P(X, Y) is generated by

 $Y = f(X) + N_Y$

with independent $X \sim \mathcal{N}(0, \sigma_X^2)$ and $N_Y \sim \mathcal{N}(0, \sigma_{N_Y}^2)$.

Then

 $\inf_{Q \in \{Q: Y \to X\}} \operatorname{KL}(P \mid\mid Q) \geq \operatorname{KL}(P(Y) \mid\mid \mathcal{N}(0, \operatorname{var} Y))$

Proposition

Assume P(X, Y) is generated by

 $Y = f(X) + N_Y$

with independent $X \sim \mathcal{N}(0, \sigma_X^2)$ and $N_Y \sim \mathcal{N}(0, \sigma_{N_Y}^2)$.

Then

 $\inf_{Q \in \{Q: Y \to X\}} \operatorname{KL}(P \mid\mid Q) \geq \operatorname{KL}(P(Y) \mid\mid \mathcal{N}(0, \operatorname{var} Y))$

- gives us finite sample guarantees
- model misspefication: how much non-Gaussianity can we allow for

Proposition

Assume P(X, Y) is generated by

 $Y = f(X) + N_Y$

with independent $X \sim \mathcal{N}(0, \sigma_X^2)$ and $N_Y \sim \mathcal{N}(0, \sigma_{N_Y}^2)$.

Then

 $\inf_{Q \in \{Q: Y \to X\}} \operatorname{KL}(P \mid \mid Q) \geq \operatorname{KL}(P(Y) \mid \mid \mathcal{N}(0, \operatorname{var} Y))$

- gives us finite sample guarantees
- model misspefication: how much non-Gaussianity can we allow for
- Question: $\inf_{Q \in \{Q: Y \to X\}} \operatorname{KL}(P || Q) = \dots$?

Can we reconstruct the whole causal network?

Theorem

let $P(X_1, ..., X_p)$ be generated by an additive noise model (+Gaussian)

$$X_i = f_i(X_{\mathbf{PA}_i}) + N_i$$

with jointly independent $N_i \sim \mathcal{N}(0, \sigma_i^2)$ and differentiable, non-linear f_i . Then we can identify the corresponding DAG from $P(X_1, \ldots, X_p)$.

JP, J. Mooij, D. Janzing and B. Schölkopf: Causal Discovery with Continuous Additive Noise Models, to appear in JMLR

Theorem

let $P(X_1, ..., X_p)$ be generated by an additive noise model (+Gaussian)

$$X_i = f_i(X_{\mathbf{PA}_i}) + N_i$$

with jointly independent $N_i \sim \mathcal{N}(0, \sigma_i^2)$ and differentiable, non-linear f_i . Then we can identify the corresponding DAG from $P(X_1, \ldots, X_p)$.

JP, J. Mooij, D. Janzing and B. Schölkopf: Causal Discovery with Continuous Additive Noise Models, to appear in JMLR

Surprising: This follows from identifiability in the bivariate case.

Theorem

let $P(X_1, ..., X_p)$ be generated by a causal additive model (+Gaussian)

$$X_i = \sum_{k \in \mathbf{PA}_i} f_{i,k}(X_k) + N_i$$

with jointly independent $N_i \sim \mathcal{N}(0, \sigma_i^2)$ and differentiable, non-linear $f_{i,k}$. Then we can identify the corresponding DAG from $P(X_1, \ldots, X_p)$.

JP, J. Mooij, D. Janzing and B. Schölkopf: Causal Discovery with Continuous Additive Noise Models, to appear in JMLR

Surprising: This follows from identifiability in the bivariate case.

Can we reconstruct the whole causal network?

Given $\hat{P}_n(X_1, \ldots, X_4)$. What now?

Can we reconstruct the whole causal network?

Given $\hat{P}_n(X_1, \ldots, X_4)$. What now?

Consider model classes

 $Q_G := \{Q : Q \text{ generated by a causal additive model CAM with DAG } G\}$

Optimize

 $\min_{\text{DAG } G} \inf_{Q \in \mathcal{Q}_G} \operatorname{KL}(\hat{P}_n || Q)$

Given $\hat{P}_n(X_1,\ldots,X_4)$. What now?

Consider model classes

 $Q_G := \{Q : Q \text{ generated by a causal additive model CAM with DAG } G\}$

Optimize

$$\min_{\text{DAG } G} \inf_{Q \in \mathcal{Q}_G} \operatorname{KL}(\hat{P}_n || Q)$$

$$\underset{\text{likelihood}}{\overset{\text{max.}}{\longrightarrow}} \min_{\text{DAG G}} \sum_{i=1}^{p} \log \operatorname{var}(\operatorname{residuals}_{\mathbf{PA}_{i}^{G} \to X_{i}})$$

Given $\hat{P}_n(X_1,\ldots,X_4)$. What now?

Consider model classes

 $Q_G := \{Q : Q \text{ generated by a causal additive model CAM with DAG } G\}$

Optimize

$$\min_{\text{DAG } G} \inf_{Q \in \mathcal{Q}_G} \operatorname{KL}(\hat{P}_n || Q)$$

$$\underset{\mathsf{DAG } G}{\overset{\mathsf{max.}}{\underset{\mathsf{DAG } G}{\overset{\mathsf{p}}{\underset{i=1}{\overset{\mathsf{p}}{\underset{j=1}{\overset{\mathsf{log } \mathsf{var}}{(\mathrm{residuals}_{\mathsf{PA}_{i}^{G} \rightarrow X_{i}})}}}}$$

Wait, there is no penalization on the number of edges! ~> fully connected graph, i.e. orderings There are

18676600744432035186664816926721

DAGs with 13 nodes.

There are

18676600744432035186664816926721

DAGs with 13 nodes. There are only :-)

6227020800

orderings.

Idea: find causal order

Find the order of variables that maximizes likelihood and then perform classical variable selection. This is consistent

There are

18676600744432035186664816926721

DAGs with 13 nodes. There are only :-)

6227020800

orderings.

Idea: find causal order

Find the order of variables that maximizes likelihood and then perform classical variable selection. This is consistent (but intractable).

P. Bühlmann, JP and J. Ernest: CAM: Causal add. models, high-dim. order search and penalized regression, submitted



STEP 1: Greedy Addition. Include the edge that leads to the largest increase of the log-likelihood.

Theorem (Some correctness of greedy search for CAM)

Assume that the skeleton of the correct DAG does not contain any cycles (plus some mild conditions/modifications).

Greedy addition then yields a correct causal order (in population).

JP, S. Balakrishnan and M. Wainwright: in progress.

STEP 1: Greedy Addition. Include the edge that leads to the largest increase of the log-likelihood.

Can we reconstruct the whole causal network?



STEP 2: Variable Selection. For each node, remove non-relevant edges.

Can we reconstruct the whole causal network?



Easy to add for high-dim data: **STEP 0:** Preliminary Neighbourhood Selection.

Can we reconstruct the whole causal network? Simulations

p = 100, n = 200, functions drawn from Gaussian Process



Can we reconstruct the whole causal network? Simulations

p = 100, n = 200, functions drawn from Gaussian Process



JP and P. Bühlmann: Structural Intervention Distance (SID) for Evaluating Causal Graphs, arXiv

Can we reconstruct the whole causal network? Simulations

p = 10, n = 200, linear functions



Can we reconstruct the whole causal network? Real Data: Arabidopsis thaliana



Wille et al.: Sparse graphical Gaussian modeling of the isoprenoid gene network in arabidopsis thaliana. Genome Biology, 5(11),

Can we reconstruct the whole causal network? Real Data: Arabidopsis thaliana





Wille et al.: Sparse graphical Gaussian modeling of the isoprenoid gene network in arabidopsis thaliana. Genome Biology, 5(11),

Jonas Peters (ETH Zurich)

Can we reconstruct the whole causal network? Real Data: A ground truth?

Assume we are given **observational** data (j = 1...5000; k = 1...120)

 O_{kj} : expression level of gene j in observation k,

Can we reconstruct the whole causal network? Real Data: A ground truth?

Assume we are given **observational** data $(j = 1 \dots 5000; k = 1 \dots 120)$

 O_{kj} : expression level of gene j in observation k,

and some **interventional** data ($j = 1 \dots 5000$; $i = \dots$ in total 400 genes)

 A_{ij} : expression level of gene *j* under a knock-down of gene *i*.

Can we reconstruct the whole causal network? Real Data: A ground truth?

Assume we are given **observational** data $(j = 1 \dots 5000; k = 1 \dots 120)$

 O_{kj} : expression level of gene j in observation k,

and some **interventional** data $(j = 1 \dots 5000; i = \dots$ in total 400 genes)

 A_{ij} : expression level of gene *j* under a knock-down of gene *i*.

How do we evaluate causal inference methods?

Idea: Additive Noise Models

Structural assumptions like additive noise models lead to identifiability:

$$X_i = f_i(X_{pa(i)}) + N_i$$

Idea: causal order + variable selection

Find the order of variables that maximizes likelihood (by greedily adding edges) and then perform classical variable selection.

Idea: Additive Noise Models

Structural assumptions like additive noise models lead to identifiability:

$$X_i = f_i(X_{pa(i)}) + N_i$$

Idea: causal order + variable selection

Find the order of variables that maximizes likelihood (by greedily adding edges) and then perform classical variable selection.

Open Questions

- Identifiability (e.g. in terms of KL-distances) depending on non-linearity, for example.
- Hidden variables
- Do the assumptions (roughly) hold in practice?
- How do we evaluate causal inference methods?

Idea: Additive Noise Models

Structural assumptions like additive noise models lead to identifiability:

$$X_i = f_i(X_{pa(i)}) + N_i$$

Idea: causal order + variable selection

Find the order of variables that maximizes likelihood (by greedily adding edges) and then perform classical variable selection.

Open Questions

- Identifiability (e.g. in terms of KL-distances) depending on non-linearity, for example.
- Hidden variables
- Do the assumptions (roughly) hold in practice?
- How do we evaluate causal inference methods?

Dankeschön!!

Jonas Peters (ETH Zurich)



7th June 2014