Session 5: Learning Bayesian Networks and Causal Discovery

Marek J. Drużdżel

Wydział Informatyki

Politechnika Białostocka

<u>m.druzdzel@pb.edu.pl</u> http://aragorn.wi.pb.bialystok.pl/~druzdzel/

1)

Course schedule

Day 2

50

Session 1: Introduction to probabilistic (Bayesian) modeling and inference
 Session 2: Bayesian networks
 Session 3: Building Bayesian networks
 Session 4: Hands-on exercises (Bayesian networks)

Session 5: Learning: structure/causal discovery, parameter learning, model validation techniques
Session 6: Hands-on exercises (learning)
Session 7: Decision analysis: expected utility theory, utility elicitation, influence diagrams
Session 8: Hands-on exercises (decision analysis)

Session overview

- Motivation
- Causality and probability
- Constraint-based learning
- Bayesian learning
- Example

φ

- Software demo
- Concluding remarks

What I want you to know after this session

- Understand the relationship between probability and causality
- Understand the principles behind learning causal models
- Be able to learn a model from data using GeNle

D

Learning Bayesian networks from data

c

There exist algorithms with a capability to analyze data, discover causal patterns in them, and build models based on these data.



Learning Bayesian Networks and Causal Discovery

Motivation

Example Software demo

Constraint-based learning Bayesian learning

Concluding remarks





Given (1), is (2) really surprising?

Causality and probability

[1]

The only reference to causality in a typical statistics textbook is: "correlation does not mean causation"

(if the textbook contains the word "causality" at all $\textcircled{\columnwidth{\odot}}$).

Many confusing substitute terms: "confounding factor," "latent variable," "intervening variable," etc.

What does correlation mean then (with respect to causality)?

The goal of experimental design is often to establish (or disprove) causation. We use statistics to interpret the results of experiments (i.e., to decide whether a manipulation of the independent variable caused a change in the dependent variable).

How are causality and probability actually related and what does one tell us about the other?

Not knowing this constitutes a handicap!

Causality and probability

Causality and probability are closely related and their relation should be made clear in statistics.

Probabilistic dependence is considered a necessary condition for establishing causation (is it sufficient?).



50

Flu and fever are correlated because flu may cause fever. A cause can cause an effect but it does not have to. Causal connections result in probabilistic dependencies (or correlations in linear case).

Motivation Constraint-based learning Bayesian learning Example Software demo Concluding remarks

Causal graphs

Acyclic directed graphs (hence, no time and no dynamic reasoning) representing a snapshot of the world at a given time. Nodes are random variables and arcs are direct causal dependencies between them.

Causal connections result in *correlation* (in general *probabilistic dependence*).

- glass on the road will be correlated with flat tire
- glass on the road will be correlated with noise
- bumpy feeling will be correlated with noise

ф



Causal Markov condition

56

An axiomatic condition describing the relationship between causality and probability.

A variable in a causal graph is probabilistically independent of its non-descendants given its immediate predecessors.

Axiomatic, but used by almost everybody in practice and no convincing counter examples to it have been shown so far (at least outside the quantum world).

Markov condition: Implications

Variables A and B are probabilistically dependent if there exists a directed active path from A to B or from B to A: Thorns on the road are correlated with car damage because there is a directed path from thorns to car damage.

ယ



Markov condition: Implications

Variables A and B are probabilistically dependent if there exists a C such that there exists a directed active path from C to A and there exists a directed active path from C to B: Car damage is correlated with noise because there is a directed path from flat tire to both (flat tire is a common cause of both).

1



Motivation Constraint-based learning Bayesian learning Example Software demo Concluding remarks

Markov condition: Implications

Variables A and B are probabilistically dependent if there exists a D such that D is observed (conditioned upon) and there exists a C such that A is dependent on C and there exists a directed active path from C to D and there exists an E such that B is dependent on E and there exists a directed active path from E to D: Nails on the road are correlated with glass on the road given flat tire because there is a directed path from glass on the road to flat tire and from nails on the road to flat tire and flat tire is observed (conditioned upon).

ф



Markov condition: Summary of implications

Variables A and B are probabilistically dependent if:

- there exists a directed active path from A to B or there exists a directed active path from B to A
- there exists a C such that there exists a directed active path from C to A and there exists a directed active path from C to B
- there exists a D such that D is observed (conditioned upon) and there exists a C such that A is dependent on C and there exists a directed active path from C to D and there exists an E such that B is dependent on E and there exists a directed active path from E to D

Motivation Constraint-based learning Bayesian learning Example Software demo Concluding remarks

Markov condition: Conditional independence

Once we know all direct causes of an event E, the causes and effects of those causes do not tell anything new about E and its successors.

(also known as "screening off")

E.g.,

b

- Glass and thorns on the road are independent of noise, bumpy feeling, and steering problems conditioned on flat tire.
- Noise, bumpy feeling, and steering problems become independent conditioned on flat tire.



Intervention

10

Manipulation theorem [Spirtes, Glymour & Scheines 1993]:

Given an external intervention on a variable A in a causal graph, we can derive the posterior probability distribution over the entire graph by simply modifying the conditional probability distribution of A.

If this intervention is strong enough to set A to a specific value, we can view this intervention as the only cause of A and reflect this by removing all edges that are coming into A. Nothing else in the graph needs to be modified.



Motivation Constraint-based learning Bayesian learning Example Software demo Concluding remarks

KER-

A cancer tumor!

No?! Yes!!

Wooaah!

Suicide eliminates cancer as a cause of this brave samurai's death.

1)

Intervention: Example

Intervention: Example

φ

Making the tire flat with a knife makes glass, thorns, nails, and what-haveyou irrelevant to flat tire. The knife is the only cause of flat tire.



Selection bias

(1)

Observing correlation is in general not enough to establish causality.



- If we do not randomize, we run the danger that there are common causes between smoking and lung cancer (for example genetic factors).
- These common causes will make smoking and lung cancer dependent.
- It may, in fact, also be the case that lung cancer causes smoking.
- This will also make them dependent without smoking causing lung cancer.

Experimentation

Empirical research is usually concerned with testing causal hypotheses.

Smoking and lung cancer are correlated.

Can we reduce the incidence of lung cancer by reducing smoking? In other words: Is smoking a cause of lung cancer?

Each of the following causal structures is compatible with the observed correlation:





- In a randomized experiment, coin becomes the only cause of smoking.
- Smoking and lung cancer will be dependent only if there is a causal influence from smoking to lung cancer.
- If $Pr(C|S) \neq Pr(C|\sim S)$ then smoking is a cause of lung cancer.

q

• Asbestos will simply cause variability in lung cancer (add noise to the observations).

But, can we really experiment in this domain?

Motivation Constraint-based learning Bayesian learning Example Software demo Concluding remarks

Science by observation

"... Correlation between smoking and lung cancer means as much as correlation between apple imports and raise of divorce ..."



Sir Ronald A. Fisher, a prominent statistician, father of experimental design



"... George Bush taking credit for the end of the cold war is like a rooster taking credit for the daybreak ..."



Vice-president Al Gore towards vice –president Dan Quayle during their first (vice) presidential debate, Fall 1992

Science by observation

(1)

- Experimentation is not always possible.
- We can do quite a lot by just observing.
- Assumptions are crucial in both experimentation and observation, although they are usually stronger in the latter.
- New methods in causal discovery: squeezing data to the limits

Search the data for independence relations to give us a clue about the causal relations [Spirtes, Glymour, Scheines 1993].

Bayesian search learning

56

Search over the space of models and score each model using the posterior probability of the model given the data [Cooper & Herskovitz 1992; many others].

(p

"Correlation does not imply causation"

True but only in limited settings (e.g., two variables) and typically abused by authors of college textbooks ©.

If x and y are dependent, we can indeed simplify the causal picture to four simplified cases:



Motivation Constraint-based learning Bayesian learning Example Software demo Concluding remarks

Constraint search-based learning

Not necessarily true in case of three variables:

x and z are dependent y and z are dependent x and y are independent x and y are dependent given z



Foundations of constrain-based search causal discovery

Markov Condition:
 ∧ structure ⇒ independence in data.

The causal graph determines what is independent.

55

All independences in the data are structural, i.e., are consequences of Markov condition.



Violations of faithfulness condition

ф

Faithfulness assumption is more controversial. While every scientist makes it in practice, it does not need to hold.



Given that HIV virus infection has not taken place, needle sharing is independent from intercourse.



Exam Performance

Good 50% Poor 50%

φ

The effect of staying up late before the exam on the exam performance may happen to be zero: being tired may cancel out the effect of more knowledge. But is it likely?

All possible networks ...

φ



... can be divided into equivalence classes

Theorems useful in search

Theorem 1 (skeleton)

There is no edge between X and Y if and only if X and Y are independent given *any* subset (including the null set) of the other variables.

Theorem 2 (v-structures)

11

If X—Y — Z, X and Z are not adjacent, and X and Z are independent given some set W, then $X \rightarrow Y \leftarrow Z$ if and only if W does *not* contain Y.

Causal model search

10

- **1. Find (conditional) independencies in the data.**
- 2. Infer from these independencies which (classes of) causal structures could have given rise to these independencies (e.g., the PC algorithm).

PC algorithm (sketch)

Step 0:

U

Begin with a complete undirected graph.

Step 1 (Find adjacencies):

For each pair of variables <X,Y> if X and Y are independent given some subset of the other variables, remove the X–Y edge.

Step 2: (Find v-structures):

For each triple X–Y–Z, with no edge between X and Z, if X and Z are independent given some set not containing Y, then orient X–Y–Z as $X \rightarrow Y \leftarrow Z$.

Step 3 (Avoid new v-structures and cycles):

- if $X \rightarrow Y$ —Z, but there is no edge between X and Z, then orient Y–Z as Y→Z.
- if X—Z, and there is already a directed path from X to Z, then orient X Z as $X \rightarrow Z$.



PC algorithm: Example



Independencies entailed by the Markov condition:

 $\mathbf{A} \perp \mathbf{B}$ $\mathbf{A} \perp \mathbf{D} \mid \mathbf{B}, \mathbf{C}$

(0) Begin with

b



(1) From $A \perp B$, remove A - B





(3) Avoid a new v-structure $(A \rightarrow C \leftarrow D)$, Orient C –D as C \rightarrow D.





 $\begin{array}{c} A \\ \hline \\ C \\ \hline \\ B \end{array} \xrightarrow{} D \end{array}$

Patterns: Output of the PC algorithm

1

PC algorithm outputs a 'pattern', a kind of graph containing directed (\rightarrow), bi-directional (\leftrightarrow), and undirected (—) edges which represents a Markov equivalence class of Models

- A directed edge A→B in the 'pattern' indicates that there is an edge oriented A→B in every graph in the Markov equivalence class
- A bi-directional edge A↔B in the 'pattern' indicates that there is an edge between A and B in every graph in the Markov equivalence class, although its direction is impossible to establish based on the data
- An undirected edge A—B in the 'pattern', indicates that there is an edge between A and B in every graph in the Markov equivalence class, although its direction is impossible to establish based on the data; there is a possible common cause between these variables in every graph in the Markov equivalence class

Motivation Constraint-based learning Bayesian learning Example Software demo Concluding remarks

Dealing with errors in independence tests: Concluding remark Search with a varying value of statistical significance

- Independence tests performed in the first phase of the algorithm may result in Type I and Type II errors.
- It is a good practice to vary the level of statistical significance α , from very low to very high values.

56

- Graphs found with low values of α will be sparse. One can trust existence of arcs (low value of α , hard to reject null hypothesis H₀ that variables are independent; when H₀ still gets rejected, it means that the dependence was strong/robust).
- Graphs found with high values of α will be dense. One can trust absence of arcs (high value of α , easy to reject H₀ that variables are independent; when H₀ still does not get rejected, it means that the independence was strong/robust).

Continuous data

55

- Causal discovery is independent of the actual distribution of the data.
- The only thing that we need is a test of (conditional) independence.
- No problem with discrete data.
- In continuous case, we have a test of (conditional) independence (partial correlation test) when the data comes from multi-variate Normal distribution.
- Need to make the assumption that the data is multi-variate Normal.
- The discovery algorithm turns out to be very robust to this assumption [Voortman & Druzdzel, 2008].



ф





Multi-variate normality is equivalent to two conditions: (1) Normal marginals and (2) linear relationships





Linearity

ф

Multi-variate normality is equivalent to two conditions: (1) Normal marginals and (2) linear relationships

Bayesian search learning

ф

Elements of a search procedure

- A representation for the current state (a network structure.)
- A scoring function for each state (the posterior probability).
- A set of search operators.
 - AddArc(X,Y)
 - DelArc(X,Y)
 - RevArc(X,Y)
- A search heuristic (e.g., greedy search).
- The size of the search space for n variables is almost 3 to the power of Cⁿ₂ possible graphs! (e.g., for 10 variables, we have 3⁴⁵ possible graphs)

Motivation Constraint-based learning Bayesian learning Example Software demo Concluding remarks Posterior probability score

$$P(S \mid D) = \frac{P(D \mid S)P(S)}{P(D)} \propto P(D \mid S)P(S)$$

"Marginal likelihood" P(D|S):

• Given a database

ф

Assuming Dirichlet priors over parameters

$$P(D \mid S) = \prod_{i=1}^{n} \prod_{j=1}^{q_i} \frac{\Gamma(\alpha_{ij})}{\Gamma(\alpha_{ij} + N_{ij})} \cdot \prod_{k=1}^{r_i} \frac{\Gamma(\alpha_{ijk} + N_{ijk})}{\Gamma(\alpha_{ijk})}$$



Starting from a variety of different points (in this case, a variety of different graphs) increases the probability of finding the graph with a maximum score.

φ

Motivation Constraint-based learning Bayesian learning Example Software demo Concluding remarks

Constraint-based learning: Open problems

Pros:

- Efficient, O(n²) for sparse graphs.
- Hidden variables can be discovered in a modest way.
- "Older" technology, many researchers do not seem to be aware of it.

11



- Discrete independence tests are computationally intensive
 - \Rightarrow heuristic independence tests?
- Missing data is difficult to deal with
 - \Rightarrow Bayesian independence test?

Bayesian learning: Open problems

Pros:

- Missing data and hidden variables are easy to deal with (in principle).
- More flexible means of specifying prior knowledge.
- Many open research questions!

D

Cons:

- Essentially intractable.
- Search heuristics (most efficient) typically lead to local maxima.
- Monte-Carlo techniques (more accurate) are very slow for most interesting problems.

Example application

10

- Student retention in US colleges.
- Large problem for US colleges.
- Correctly predicted that the main causal factor in low student retention is the quality of incoming students.

[Druzdzel & Glymour, 1994]

Example: What causes low student retention?

• Some US colleges lose over 80% of their incoming (undergraduate) students within the first year.

φ

• Below a histogram of the 1994 retention rates of 170 US national colleges.



Everything seems to be correlated with everything. What would you suggest causes low student retention?

	spend	apret	top10	rejr	tstsc	pacc	strat	salar
spend	-							
apret	0.60 <mark>123</mark> 1	-						
top10	0.67 <mark>565</mark> 6	0.64 <mark>246</mark> 4	-					
rejr	0.63 <mark>354</mark> 4	0.51 <mark>49</mark> 58	0.64 <mark>316</mark> 3	-				
tstsc	0.71 <mark>491</mark>	0.78 <mark>218</mark> 3	0.79 <mark>880</mark> 7	0.62 <mark>860</mark> 1	-			
расс	-0. <mark>2</mark> 3673	-0.3 <mark>0</mark> 2834	-0.2 <mark>0</mark> 7505	-0.07 <mark>15207</mark>	-0.1 <mark>6</mark> 4223	-		
strat	-0 <mark>.56</mark> 1755	-0. <mark>45</mark> 8311	-0.2 <mark>4</mark> 7857	-0.2 <mark>8</mark> 3617	-0. <mark>46</mark> 5226	0.13 <mark>1</mark> 858	-	
salar	0.71 <mark>183</mark> 8	0.63 <mark>585</mark> 2	0.63 <mark>764</mark> 8	0.60 <mark>677</mark> 7	0.71 <mark>547</mark> 2	-0. <mark>3</mark> 7524	-0. <mark>34</mark> 7673	-

b

Example: What causes low student retention?

- It turns out that every model that we obtain by means of a learning procedure has a direct link between test scores and high school standing (measures of the quality of incoming students) and retention.
- This finding has been confirmed by a real-world experiment.





Motivation

Example Software demo

Constraint-based learning Bayesian learning

Concluding remarks

Learning Bayesian Networks and Causal Discovery

Some challenges

56

Scaling up -- especially Monte Carlo techniques.
Practically dealing with hidden variables -unsupervised classification.
Applying these techniques to real data and real
problems.
Hybrid techniques: Constraint-based + Bayesian
(e.g., Dash & Druzdzel, 1999).

Learning causal graphs in time-dependent domains (Dash & Druzdzel, 2002).

Learning causal graphs and causal manipulation (Dash & Druzdzel, 2002).

Learning dynamic causal graphs from time series data (Voortman, Dash & Druzdzel 2010)





The reminder of this session

ф



Concluding remarks

55

- Observation is a valid scientific method
- Observation allows often to restrict the class of possible causal structures that could have generated the data.
- Learning Bayesian networks/causal graphs is very exciting: It is a different and powerful way of doing science.
- There is a rich assortment of unsolved problems in causal discovery / learning Bayesian networks, both practical and theoretical.
- Learning has been an active area of my research (GeNIe, <u>https://www.bayesfusion.com/</u>, is a product of this work).

