

# Session 5: Learning Bayesian Networks and Causal Discovery

**Marek J. Druzdzel**

University of Pittsburgh  
School of Information Sciences  
and Intelligent Systems Program

[marek@sis.pitt.edu](mailto:marek@sis.pitt.edu)  
<http://www.pitt.edu/~druzdzel>

Summer School on “*Modeling and Decision Making Using Bayesian Statistics*”  
Aalto University, Department of Applied Mechanics, Marine Technology in Espoo, Finland  
4th – 8th June, 2012

## Course schedule

Day 1 Monday

Session 1: Introduction to Bayesian inference

Session 2: Bayesian networks

Session 3: Building Bayesian networks

Session 4: Hands-on exercises (Bayesian networks)

Day 2 Tuesday

Session 5: **Learning Bayesian networks and causal discovery**

Session 6: Decision theory and decision analysis

Session 7: Hands-on exercises (decision modeling)

Session 8: Hands-on exercises (learning)

## Session overview

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

# Learning Bayesian networks from data

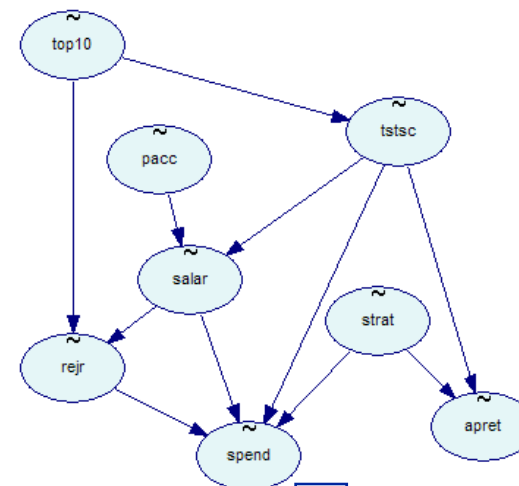
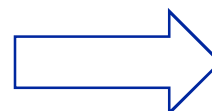
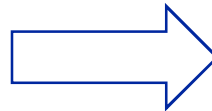
There exist algorithms with a capability to analyze data, discover causal patterns in them, and build models based on these data.

Retention.txt

	spend	apret	top10	rejr	tstsc	pacc	strat	salar
9855	52.5	15	29.474	65.063	36.887	12	60800	
10527	64.25	36	22.309	71.063	30.97	12.8	63900	
7904	37.75	26	25.853	60.75	41.985	20.3	57800	
6601	57	23	11.296	67.188	40.289	17	51200	
7251	62	17	22.635	56.25	46.78	18.1	48000	
6967	66.75	40	9.718	65.625	53.103	18	57700	
8489	70.333	20	15.444	59.875	50.46	13.5	44000	
9554	85.25	79	44.225	74.688	40.137	17.1	70100	
15287	65.25	42	26.913	70.75	28.276	14.4	71738	
7057	55.25	17	24.379	59.063	44.251	21.2	58200	
16848	77.75	48	26.69	75.938	27.187	9.2	63000	
18211	91	87	76.681	80.625	51.164	12.8	74400	
21561	69.25	58	44.702	76.25	26.689	9.2	75400	
20667	65	68	22.995	75.625	28.038	11	66200	
10684	61.75	26	8.774	66	33.99	9.5	52900	
11738	74.25	32	25.449	66.875	27.701	12	63400	
10107	74	43	11.315	71	29.096	16.2	66200	
7817	65.75	36	33.709	64.25	52.548	17.7	54600	
7050	26	11	0	55.313	55.651	18.8	59500	
9082	83.5	73	64.668	77.375	43.185	13.6	66700	
11706	60	56	16.937	73.75	39.479	12.7	62100	
7643	49.25	23	36.635	62.813	39.302	18.7	57700	
25734	90	77	67.758	80.938	44.133	10	80200	
20155	86	84	69.31	79.688	48.766	17.6	74000	
29852	94.5	84	75.009	81.313	51.363	10.6	74100	
7980	68.5	34	9.122	63.875	35.294	16.3	53100	

Row 1 of 170

data



structure

Success	0.2
Failure	0.8

	Success	Failure
Good	0.4	0.1
Moderate	0.4	0.3
Poor	0.2	0.6

numerical parameters

## The problem of learning

**Given a set of variables (a.k.a. attributes)  $X$  and a data set  $D$  of simultaneous values of variables in  $X$**

- 1. Obtain insight into causal connections among the variables  $X$  (for the purpose of understanding and prediction of the effects of manipulation)**
- 2. Learn the joint probability distribution over the variables  $X$**

## Why are we also interested in causality?

**Reason 1:** Ease of model-building and model enhancements: People think in causal terms.

**Reason 2:** Predicting the effects of manipulation.

**Given (2), (1) is not really surprising**

## Causality and probability

**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 ☺).

**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).**

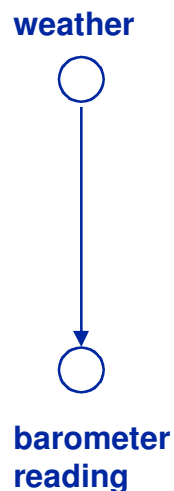
**How are causality and probability actually related and what does one tell us about the other?**

**Not knowing this constitutes a handicap!**

## Causality and probability

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

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



Weather and barometer reading are correlated **because** the weather causes the barometer reading.

A cause can cause an effect but it does not have to. Causal connections result in probabilistic dependencies (or correlations in linear case).



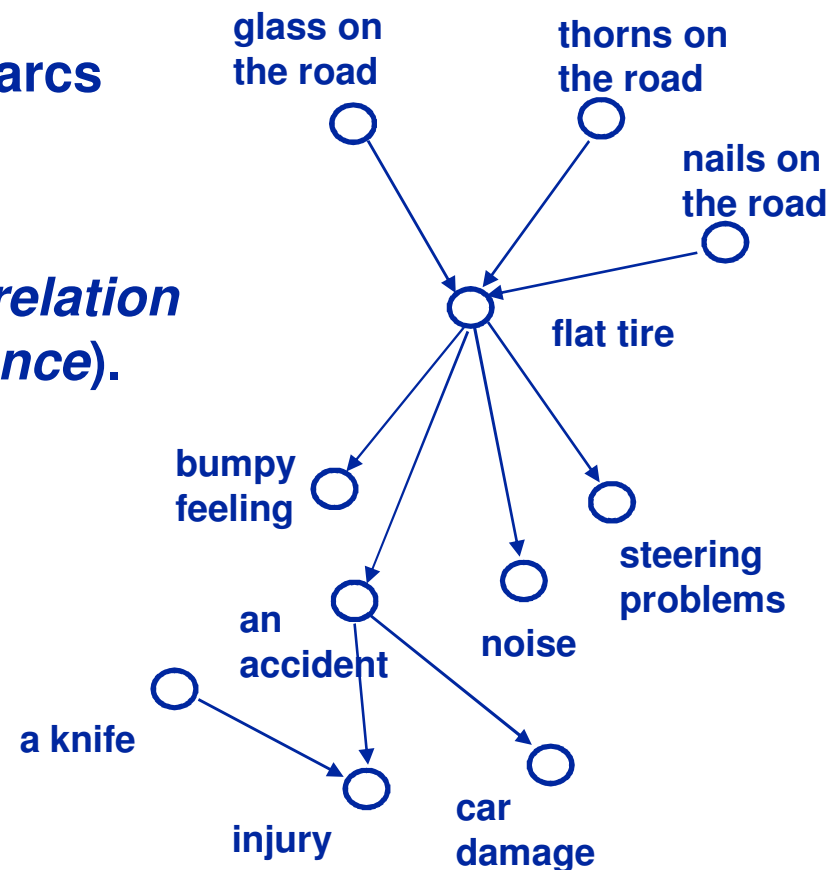
## Causal graphs

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



## Causal Markov condition

An axiomatic condition describing the relationship between causality and probability.

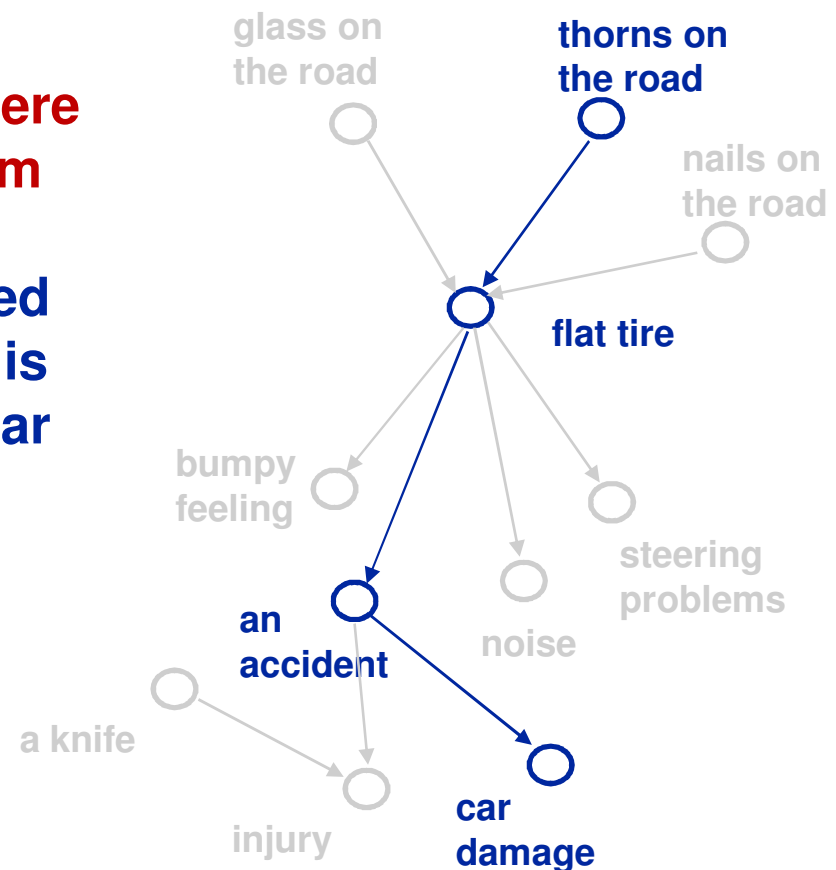
**A variable in a causal graph is probabilistically independent of its non-descendants given its immediate predecessors.**

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:**

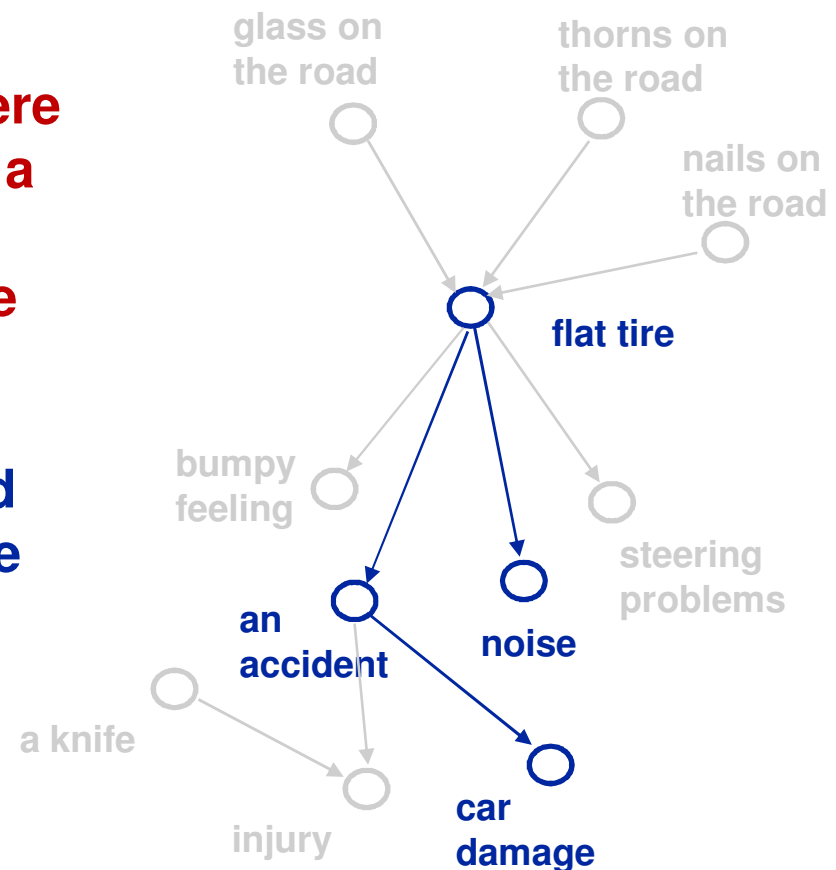
**Thorns on the road are correlated with car damage because there is a directed path from thorns to car damage.**



## 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:**

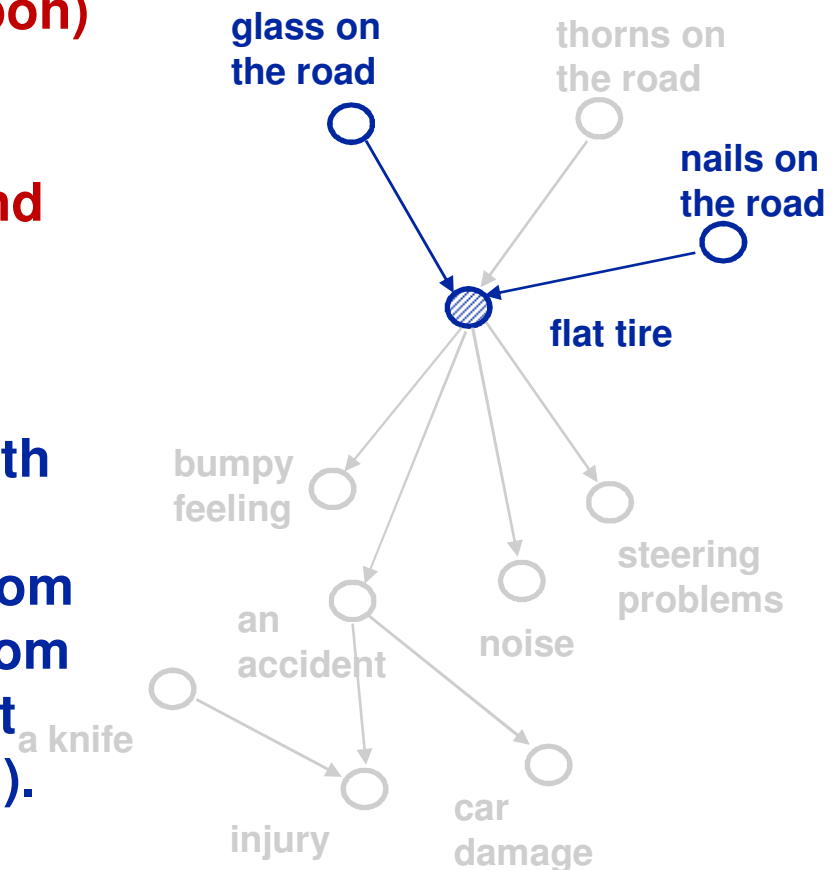
**Car damage is correlated with noise because there is a directed path from flat tire to both (flat tire is a common cause of both).**



## 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).**



## **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**

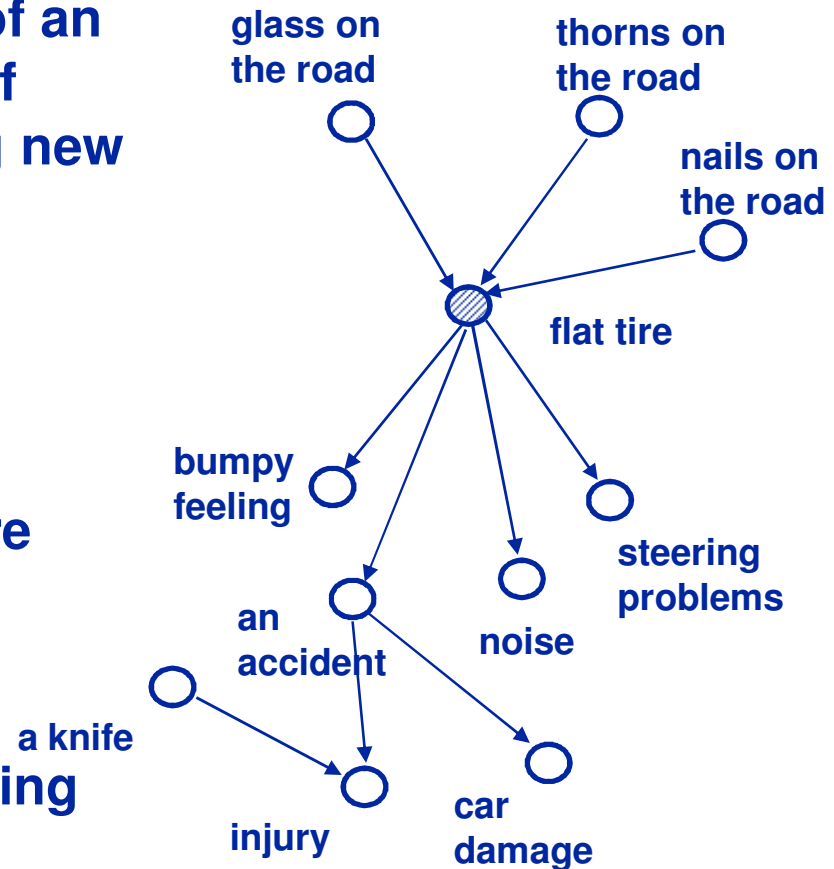
## Markov condition: Conditional independence

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.,

- 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.

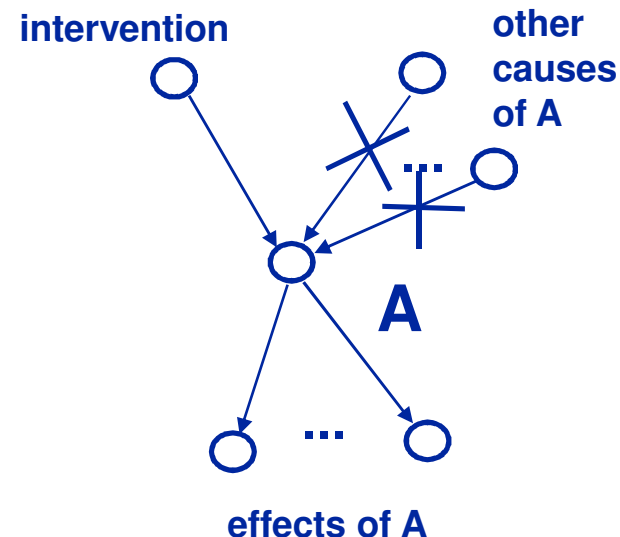


## Intervention

### 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**.

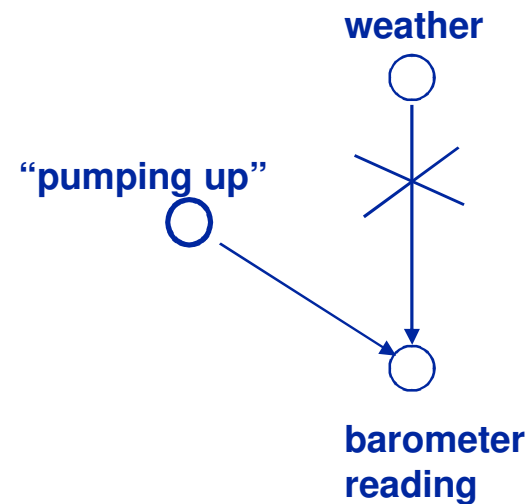
If this intervention is strong enough to set **A** to a specific value, we can view this intervention as the only cause of **A** and reflect this by removing all edges that are coming into **A**. Nothing else in the graph needs to be modified.





## Intervention: Example

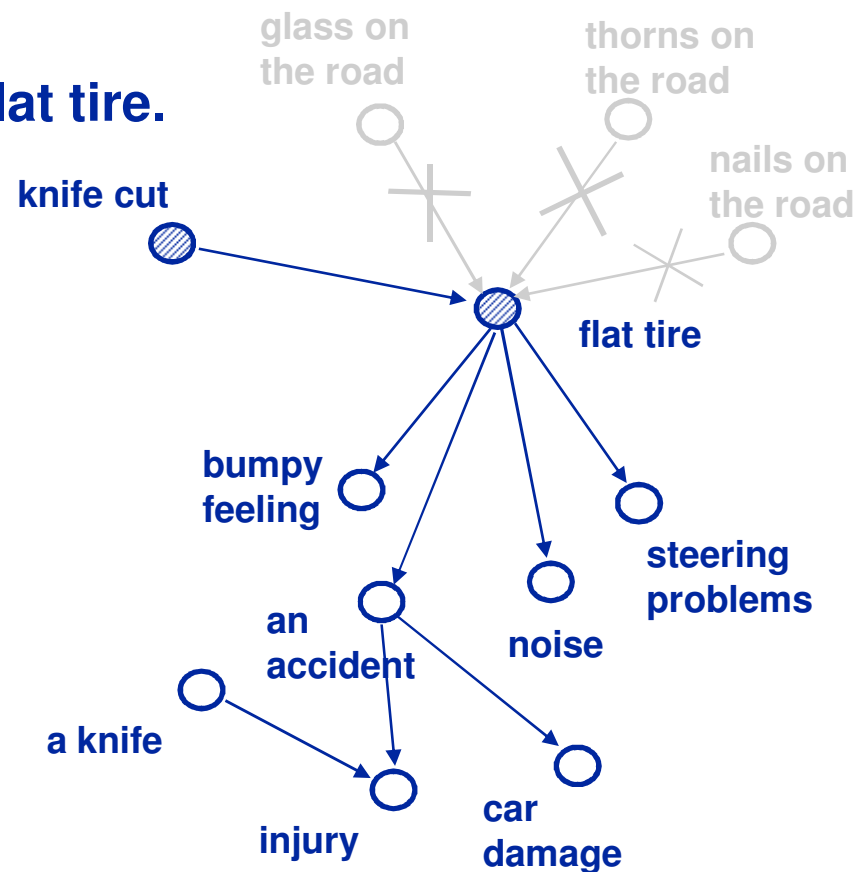
**“Pumping up” the barometer eliminates the weather as a cause of the pressure indicated by the barometer reading.**



## Intervention: Example

**Making the tire flat with a knife makes glass, thorns, nails, and what-have-you irrelevant to flat tire.**

**The knife is the only cause of flat tire.**



## Experimentation

Empirical research is usually concerned with testing causal hypotheses.

**Smoking and lung cancer are correlated.**

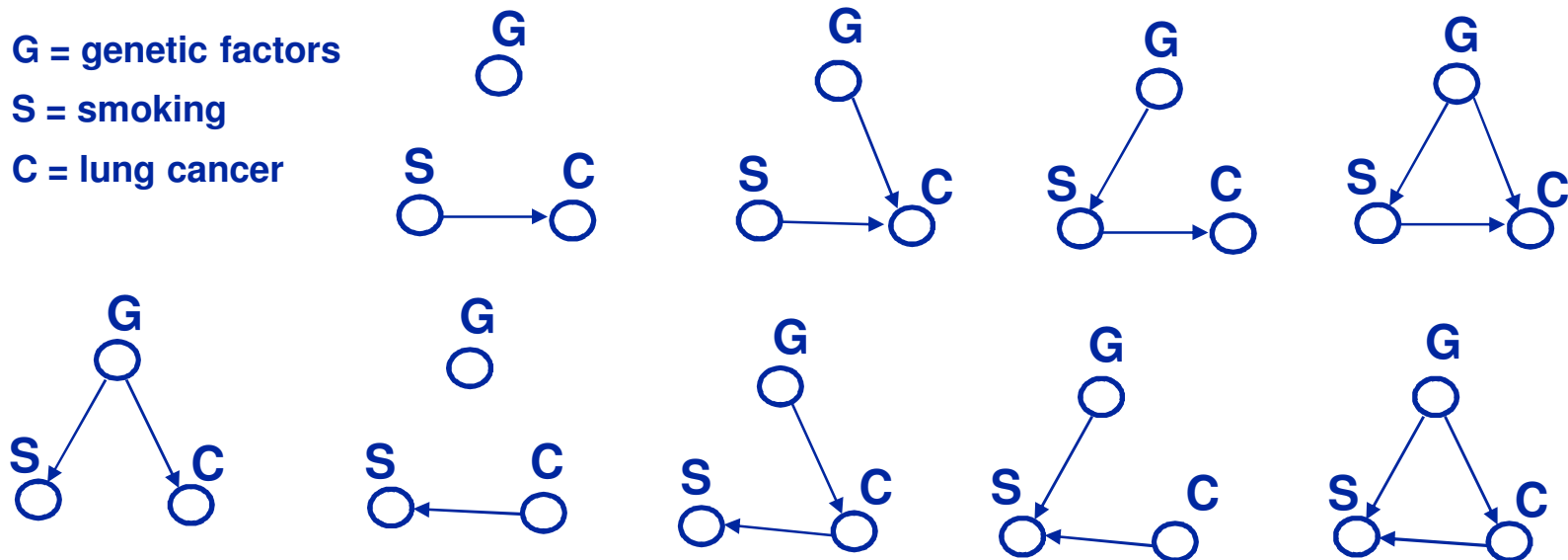
Can we reduce the incidence of lung cancer by reducing smoking?  
In other words: Is smoking **a cause** of lung cancer?

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

G = genetic factors

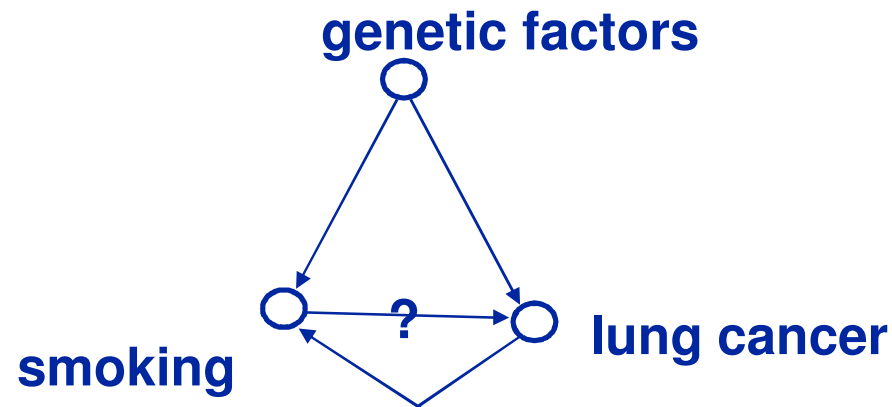
S = smoking

C = lung cancer



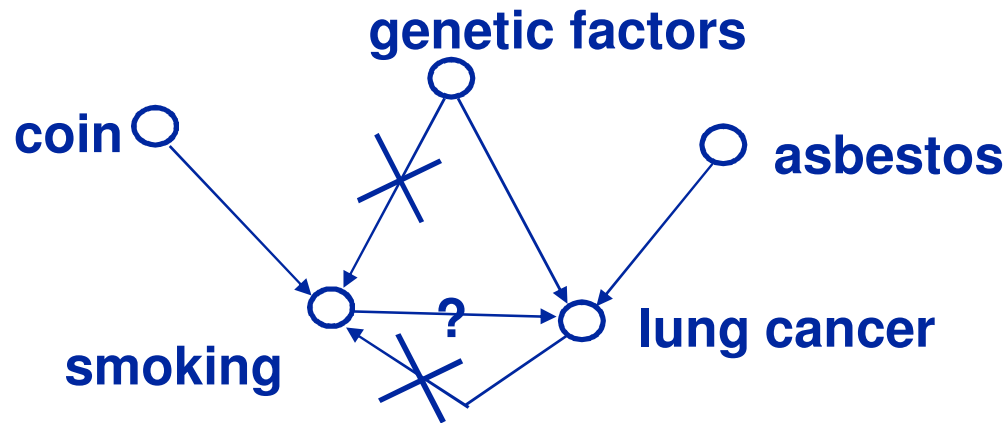
## Selection bias

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



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

**But, can we really experiment in this domain?**

## Science by observation

**“... Does smoking cause lung cancer or does lung cancer cause smoking? ...”**

**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 Dan Quayle during their first debate, Fall 1992**

- **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**

## Approaches to learning Bayesian networks

### **Constraint search-based learning**

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

### **Bayesian learning**

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].

# Constraint search-based learning



# Constraint search-based learning

## Principles:

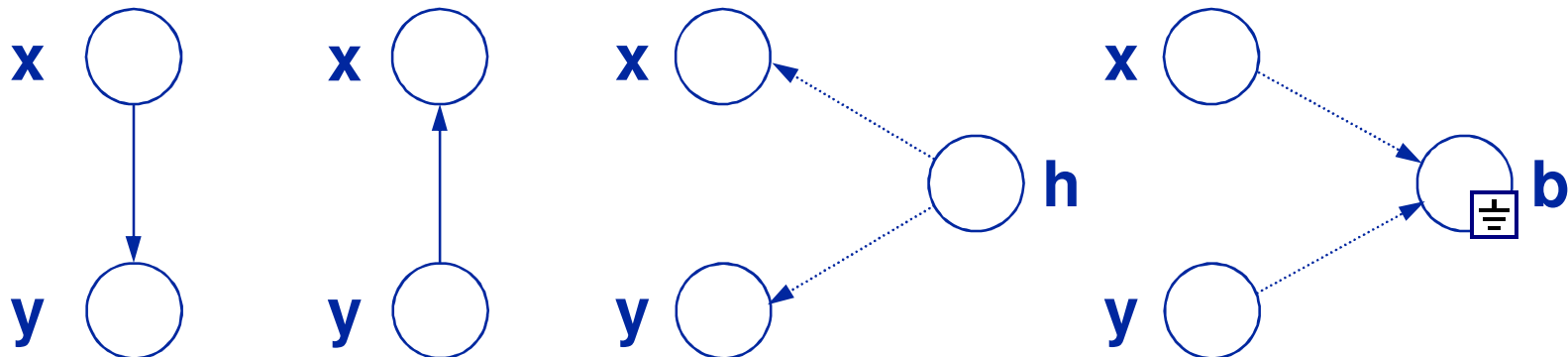
- Search for independencies among variables in the database.
- Use the *independencies* in the data to infer (lack of) *causal links* among the variables (given some basic assumptions).

## Constraint search-based learning

### “Correlation does not imply causation”

True but only in limited settings and often unfairly abused by the “statistics mafia” 😊.

If  $x$  and  $y$  are dependent, we have indeed at least four possible cases:

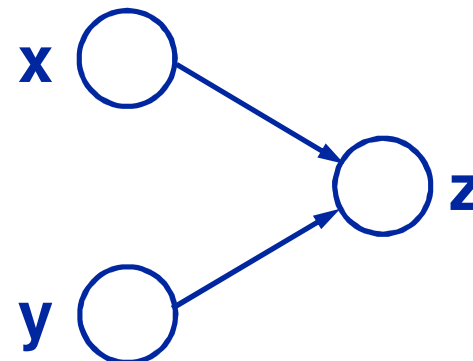


## Constraint search-based learning

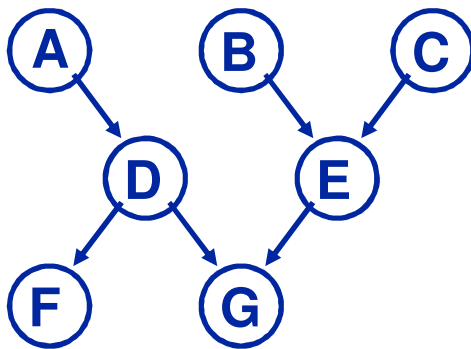
Not necessarily true in case of three variables:

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

**We can establish  
causality!**



## Foundations of causal discovery: (1) The Causal Markov Condition



Relates a causal graph to a probability distribution.

### Intuition:

In a causal graph, the parents of each node “shield” the node from their ancestors.

### Formally:

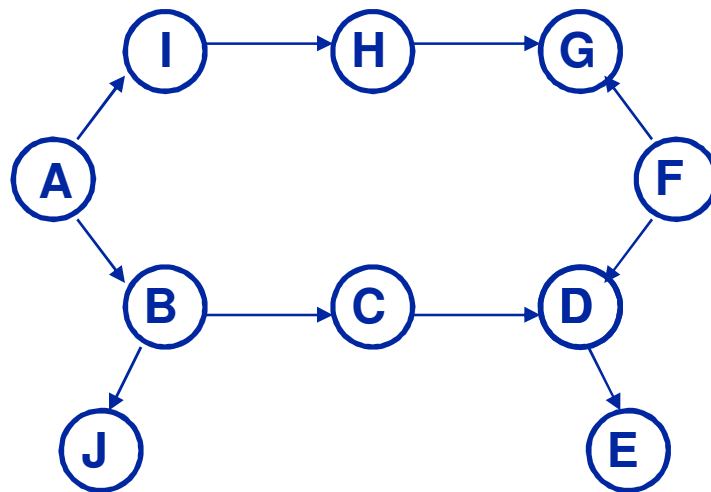
For any node  $X_i$  in the graph, we have

$$P[X_i | X', \text{Pa}(X_i)] = P[X_i | \text{Pa}(X_i)],$$

where  $\text{Pa}(X_i)$  are the parents of  $X_i$  in the graph, and  $X'$  is any set of non-descendants of  $X_i$  in the graph.

**Theorem: A causal graph obeys the Markov condition if and only if every d-separation in the graph corresponds to an independence in the probability distribution.**

## The Causal Markov Condition: d-separation



### Restatement of “the rules:”

- Each node is a “valve”
- v-structures are “off” by default
- other nodes are “on” by default
- conditioning on a node flips its state
- conditioning on a v-structure’s descendants also flips its state.

$I(B, F) ?$  **Yes**

$I(B, F \mid D) ?$  **No**

$I(B, F \mid C, D) ?$  **Yes**

## Foundations of causal discovery: (2) Faithfulness condition

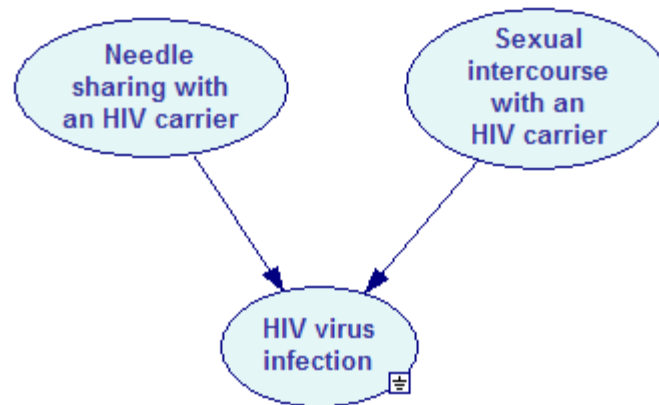
- **Markov Condition:**  
 $\text{d-separation} \Rightarrow \text{independence in data.}$
- **Faithfulness Condition:**  
 $\text{d-separation} \Leftarrow \text{independence in data.}$

**In other words:**

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

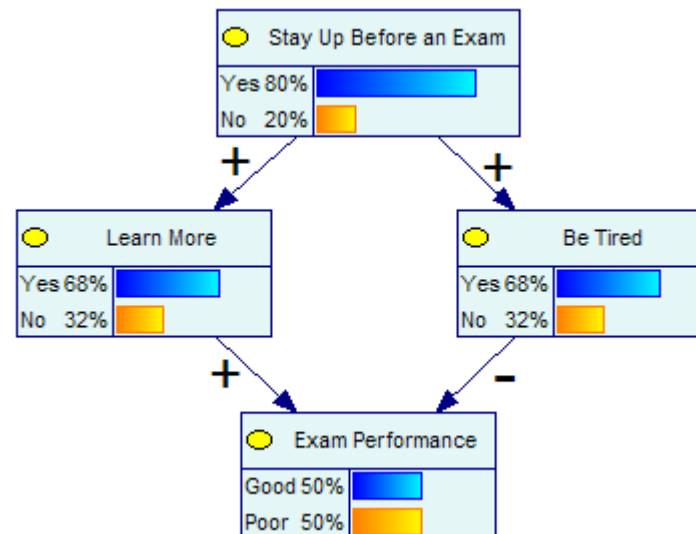
## Violations of faithfulness condition

**Faithfulness condition 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.**

## Violations of faithfulness condition

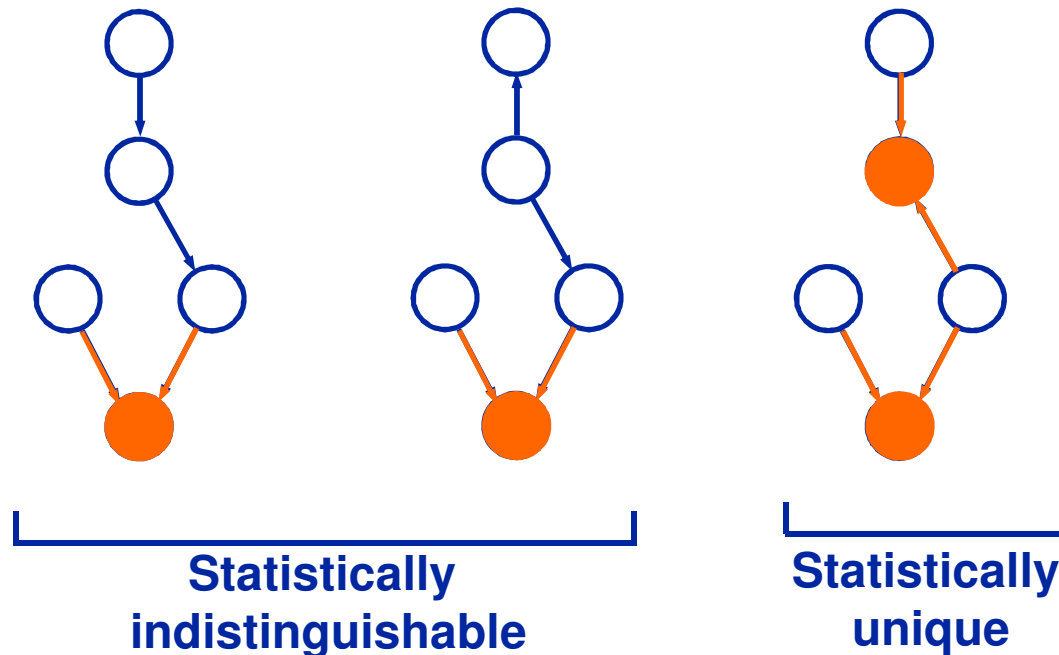


The effect of staying up late before the exam on the exam performance may happen to be zero: being tired may cancel out the effect of more knowledge. But is it likely?



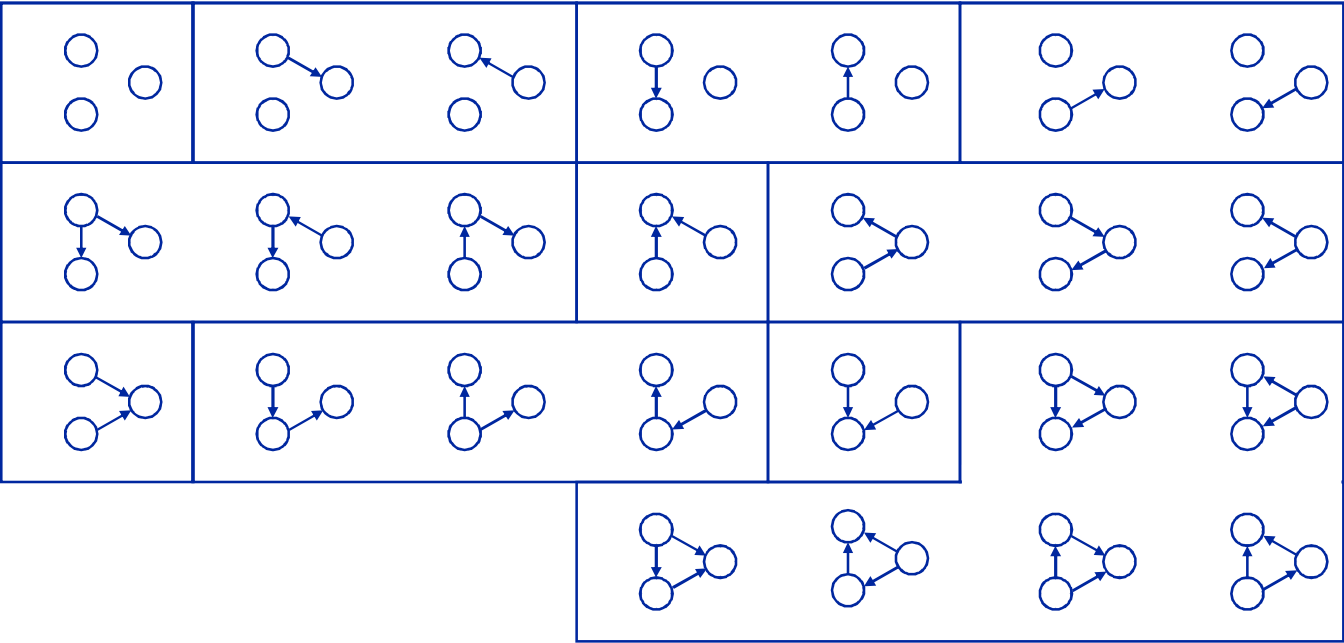
## Equivalence criterion

Two graphs are statistically indistinguishable (belong to the same equivalence class) iff they have the same adjacencies and the same “v-structures”.



# Constraint search-based learning

All possible graphs ...



... can be divided into equivalence classes

## Causal model search

1. Start with data.
2. Find conditional independencies in the data.
3. Infer which causal structures could have given rise to these independencies.

## Theorems useful in search

### Theorem 1

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

If  $X \text{---} Y \text{---} Z$ ,  $X$  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  $Y$ .

## PC algorithm

**Input:**

a set of conditional independencies

**Output:**

a “pattern” which represents a Markov equivalence class of causally sufficient causal models.

## PC algorithm (sketch)

### Step 0:

Begin with a complete undirected graph.

### Step 1 (Find adjacencies):

For each pair of variables  $\langle X, Y \rangle$  if  $X$  and  $Y$  are independent given some subset of the other variables, remove the  $X$ – $Y$  edge.

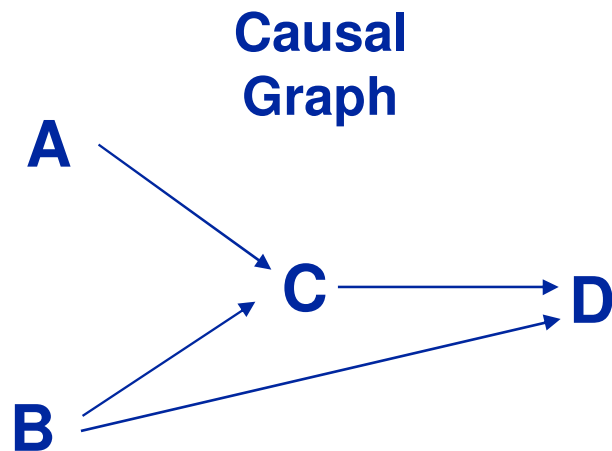
### Step 2: (Find v-structures):

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  $Y \rightarrow Z$ .
- 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

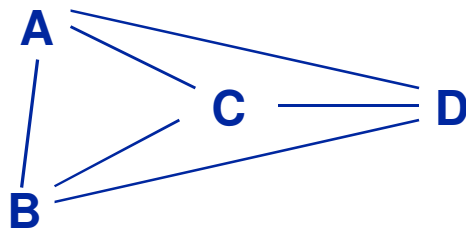


Independencies entailed by the Markov condition:

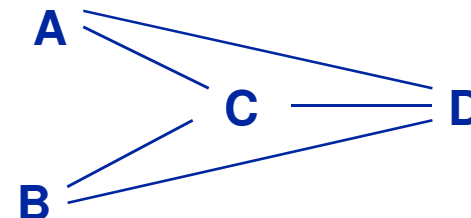
$I(A, B)$

$I(A, D \mid B, C)$

(0) Begin with

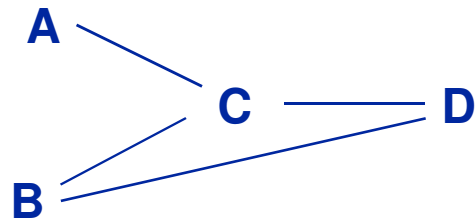


(1) From  $A \perp B$ , remove  $A-B$

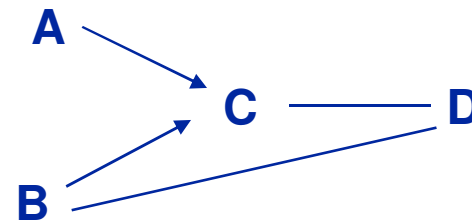


## PC algorithm: Example

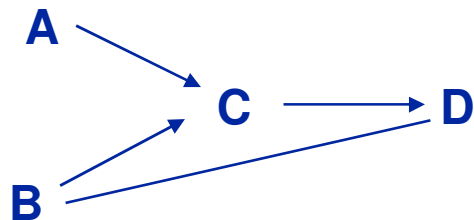
(1) From  $I(A,D \mid B,C)$ , remove  $A-D$



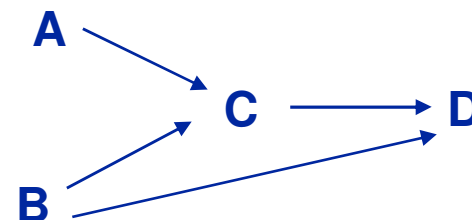
(2) From  $I(A,B)$ , orient  $A-C-B$  as  $A \rightarrow C \leftarrow B$



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



(3) Avoid a cycle  $(B \rightarrow C \rightarrow D \rightarrow B)$ ,  
Orient  $B-D$  as  $B \rightarrow D$ .





## Patterns: Output of the PC algorithm

PC algorithm outputs a ‘pattern’, a kind of graph containing directed ( $\rightarrow$ ) and undirected ( $—$ ) edges which represents a Markov equivalence class of models

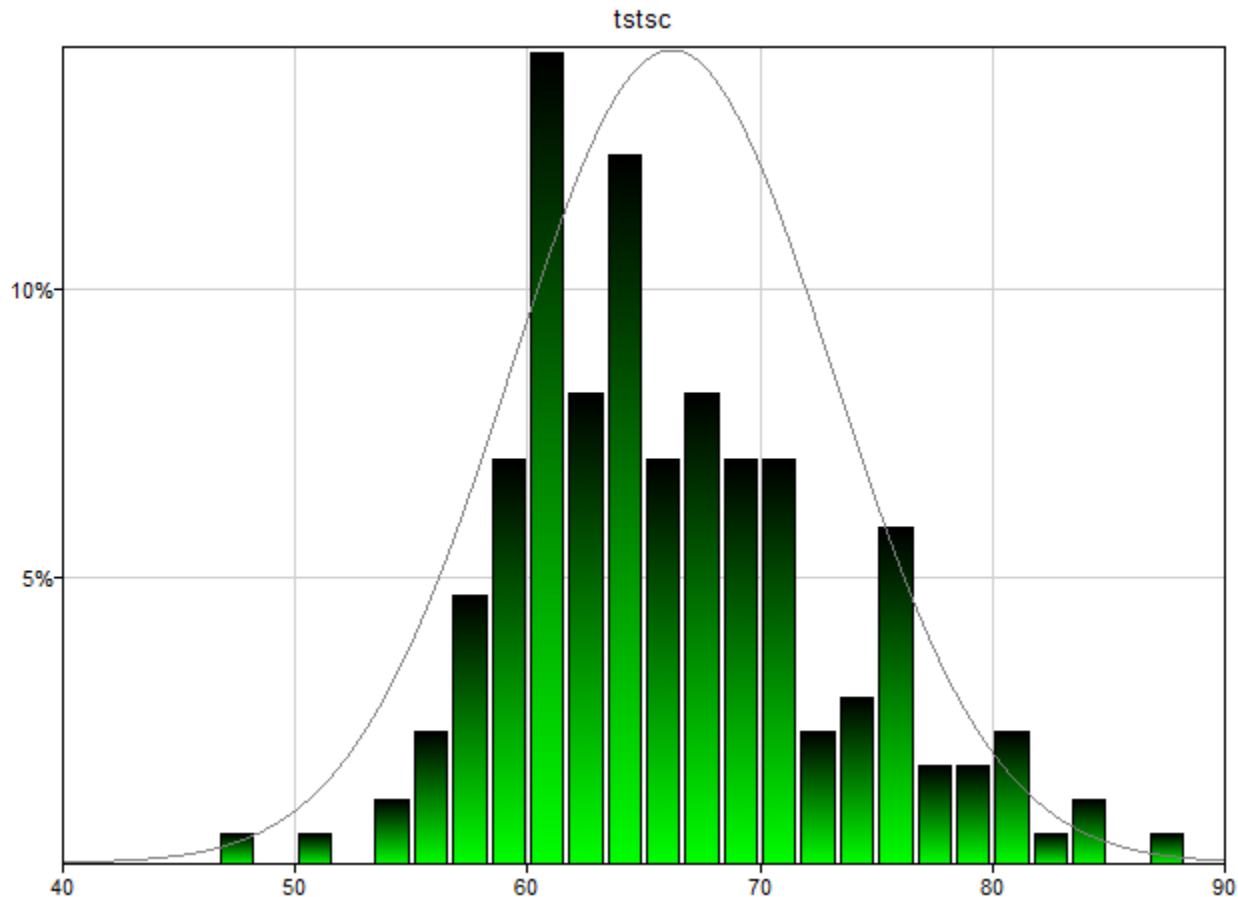
- An undirected edge  $A-B$  in the ‘pattern’, indicates that there is an edge between these variables in every graph in the Markov equivalence class
- A directed edge  $A \rightarrow B$  in the ‘pattern’ indicates that there is an edge oriented  $A \rightarrow B$  in every graph in the Markov equivalence class

## Continuous data

- 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].

# Normality

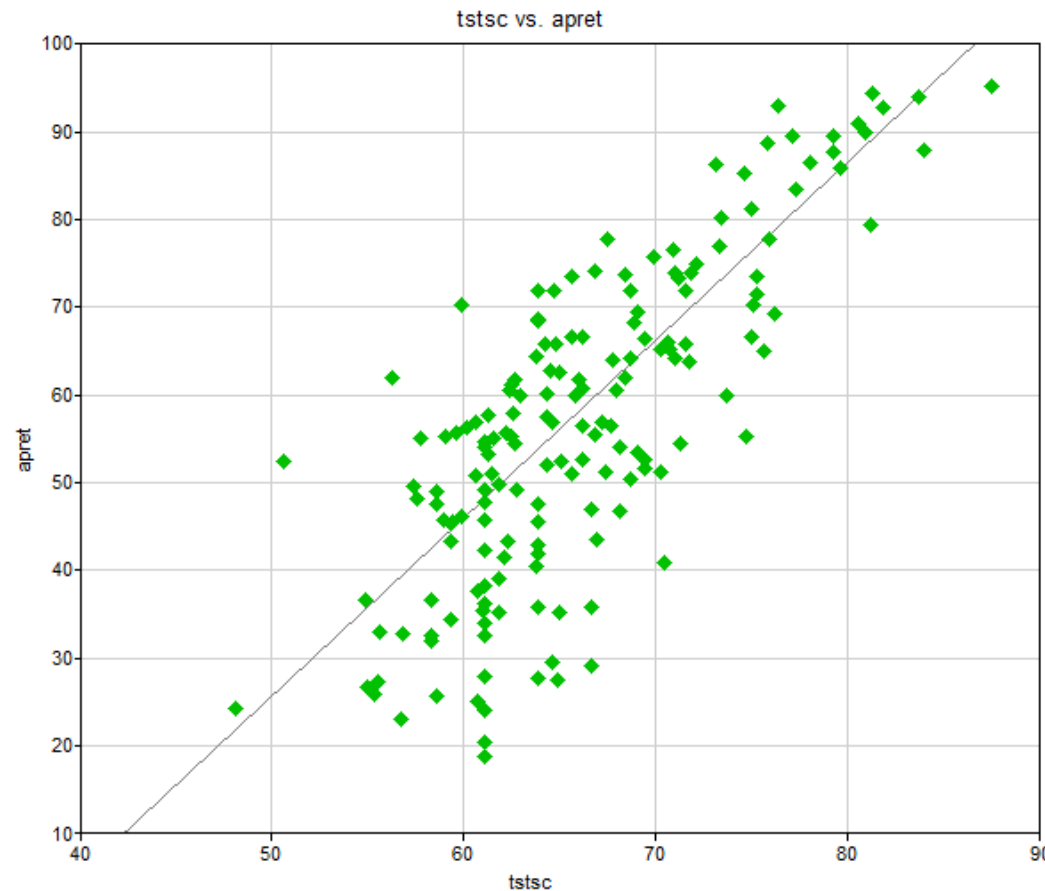
- Motivation
- Constraint-based learning
- Bayesian learning
- Example
- Software demo
- Concluding remarks



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

# Linearity

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

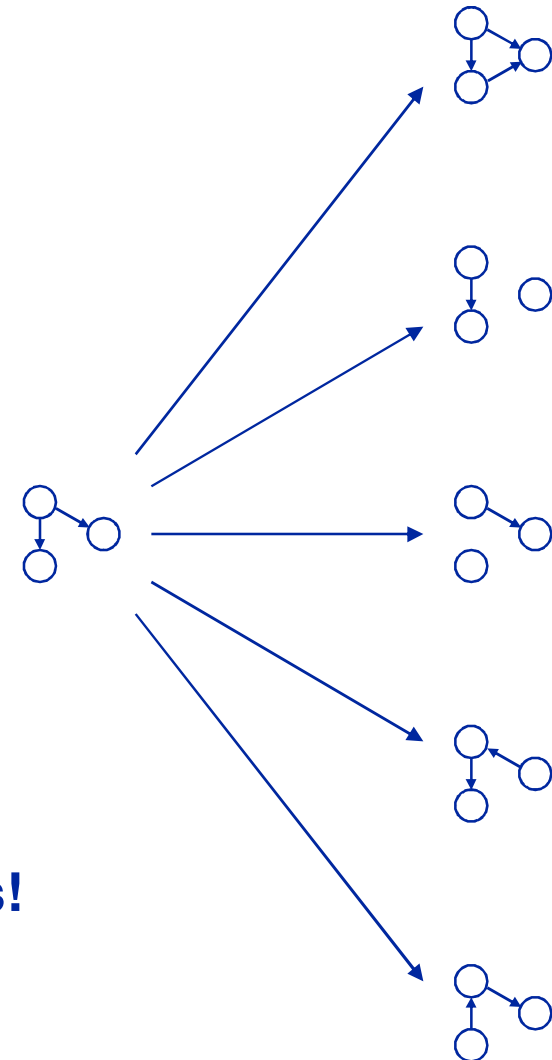


**Multi-variate normality is equivalent to two conditions:  
(1) Normal marginals and (2) linear relationships**

# Bayesian learning

## 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  $3^{C_n^2}$  possible graphs!**



## Posterior probability score

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

**“Marginal likelihood”  $P(D|S)$ :**

- Given a database
- Assuming Dirichlet priors over parameters

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

## Constraint-based learning: Open problems

### Pros:

- Efficient,  $O(n^2)$  for sparse graphs.
- Hidden variables can be discovered in a modest way.
- “Older” technology, many researchers do not seem to be aware of it.

### Cons:

- Discrete independence tests are computationally intensive  
⇒ heuristic independence tests?
- Missing data is difficult to deal with  
⇒ 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!

### Cons:

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

## Some challenges

**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).**

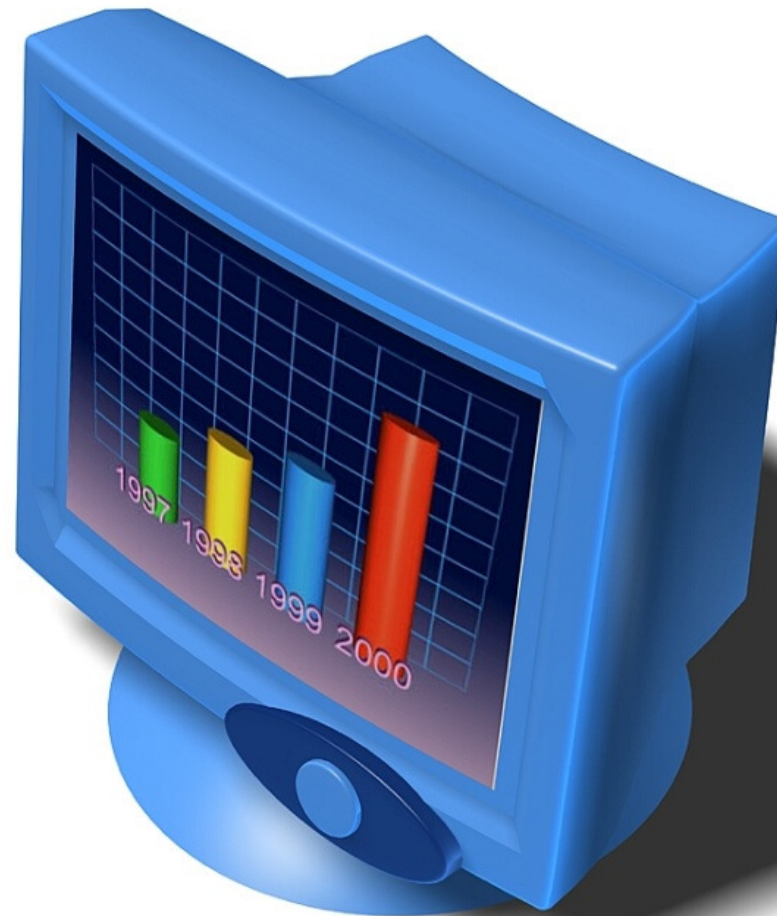
**Learning causal graphs and causal manipulation**  
**(Dash & Druzdzel, 2002).**

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



## The rest

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



## Concluding remarks

- **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.**

