

Learning Causal Bayesian Networks

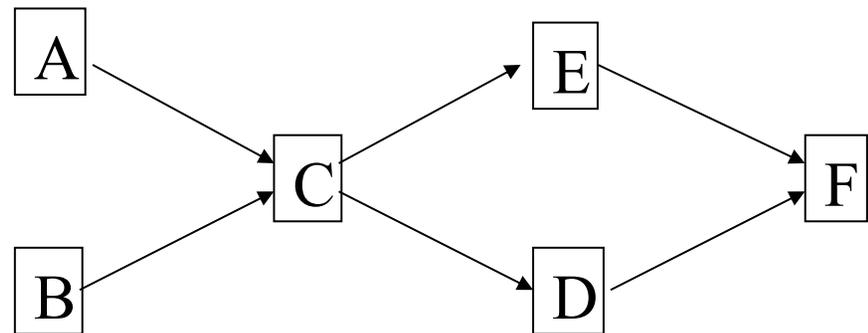
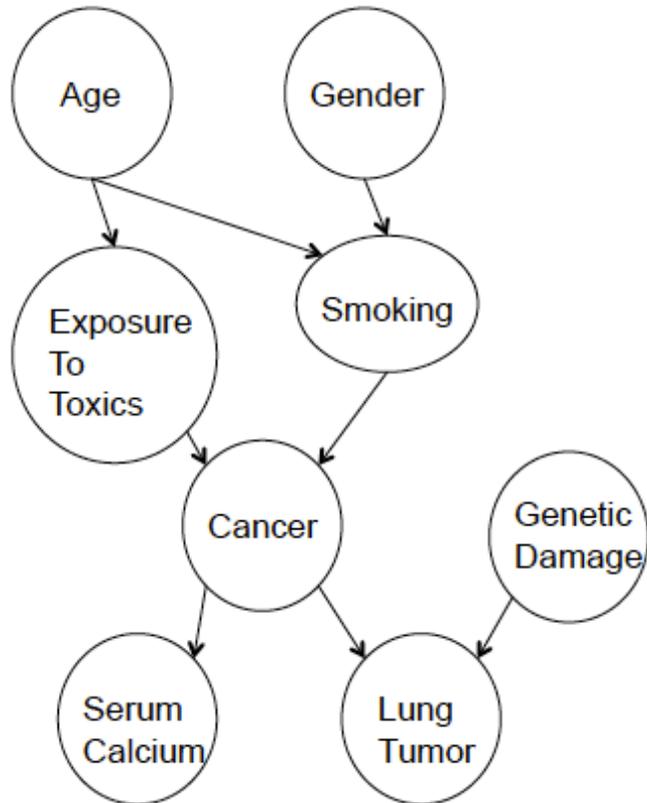
Sargur N. Srihari

University at Buffalo, The State University of New York
USA

Topics in Causality

- Causality
- Correlation vs Causality: RCT
- Importance of Causality in Transfer Learning
- Additive Noise Model
- Learning Causal Bayesian Networks

Causal Bayesian Networks



- A and B are *causally independent*;
- C, D, E, and F are *causally dependent* on A and B;
- A and B are *direct* causes of C;
- A and B are *indirect* causes of D, E and F;
- If C is prevented from changing with A and B, then A and B will no longer cause changes in D, E and F

BN is not necessarily Causal

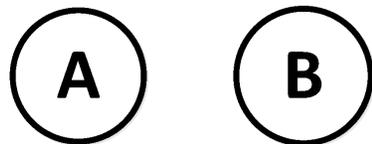
- BN is only an efficient representation of a joint distribution in terms of conditional distributions
- Several different BNs can represent the same distribution— equivalence class

Determining Causal Relationships

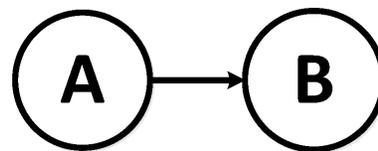
- Best way: Randomized controlled experiments
 - May be too difficult or immoral to perform
- Want to rely on “observational” data to infer causal relationships
- Current statistical methods are good at determining correlation
 - Correlation hints at causality

Relationships between events

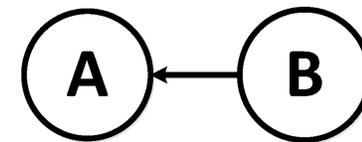
A | B



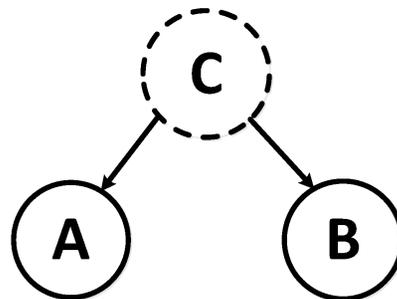
A Causes B



B Causes A



Common Causes for A and B, which do not cause each other



A = ice cream sales
B = drowning deaths
C = hot summer

Correlation is a broader concept than causation

Examples of Causal Model



- Statement 'Smoking causes cancer' implies an asymmetric relationship:
 - Smoking leads to lung cancer, but
 - Lung cancer will not cause smoking
- Arrow indicates such causal relationship
- No arrow between smoking and 'Other causes of lung cancer'
 - Means: no direct causal relationship between them

Probabilistic Causality

– Deterministic causation

- If A causes B , then A must *always* be followed by B .
 - War does not cause deaths, nor does smoking cause cancer.

– Probabilistic causation

- A probabilistically causes B if A 's occurrence increases the probability of B
 - Smoking causes cancer
- Reflects either imperfect knowledge of a deterministic system or system under study is inherently probabilistic, such as quantum mechanics

Inference with Causal Networks

1. Probabilistic Queries

– Similar to other PGMs

2. Intervention Queries

– Ideal with no other effect

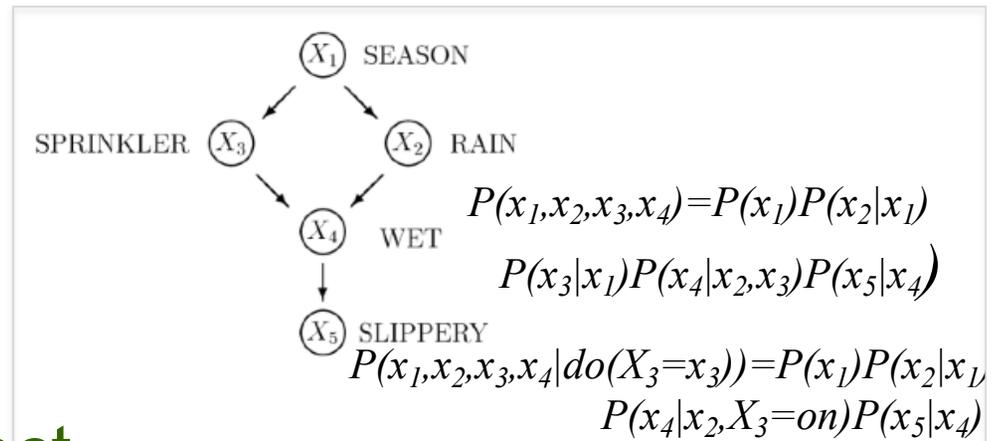
- If patient takes this medication what are chances of getting well $P(H|do(M=m^1))$

– Where H=Health. Which is different from $P(H|m^1)$

» Patients taking meds on their own are healthier

3. Contra-factual Queries

- Would the accident have happened if driver was not drunk?

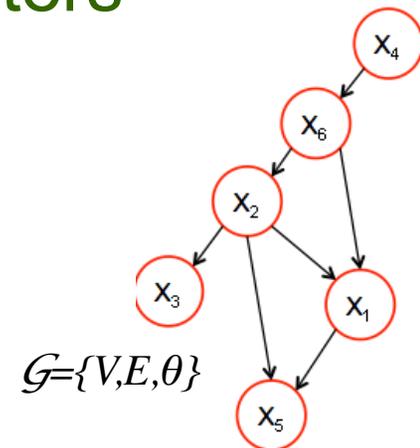


BN and Complexity of Prob Distributions

- For n variables with k states for each variable
 - Full Distribution requires $k^n - 1$ parameters
 - with $n=6, k=5$ need 15,624 parameters

R = Height Relationship of t to h	L = Shape of Loop of h	A = Shape of Arch of h	C = Height of Cross on t staff	B = Baseline of h	S = Shape of t
$r^0 = t$ shorter than h	$l^0 =$ retraced	$a^0 =$ rounded arch	$c^0 =$ upper half of staff	$b^0 =$ slanting upward	$s^0 =$ tented
$r^1 = t$ even with h	$l^1 =$ curved right side and straight left side	$a^1 =$ pointed	$c^1 =$ lower half of staff	$b^1 =$ slanting downward	$s^1 =$ single stroke
$r^2 = t$ taller than h	$l^2 =$ curved left side and straight right side	$a^2 =$ no set pattern	$c^2 =$ above staff	$b^2 =$ baseline even	$s^2 =$ looped
$r^3 =$ no set pattern	$l^3 =$ both sides curved		$c^3 =$ no fixed pattern	$b^3 =$ no set pattern	$s^3 =$ closed
	$l^4 =$ no fixed pattern				$s^4 =$ mixture of shapes

the the the



- BN Provides a factorization of joint distribution
 - Nodes are variables, edges are influences

$$P(\mathbf{x}) = P(x_4)P(x_6 | x_4)P(x_2 | x_6)P(x_3 | x_2)P(x_1 | x_2, x_6)P(x_5 | x_1, x_2)$$

θ : $4 + (3 * 24) + (2 * 125) = 326$ parameters

Organized as Six CPTs, e.g.

$$P(X_5 | X_1, X_2)$$

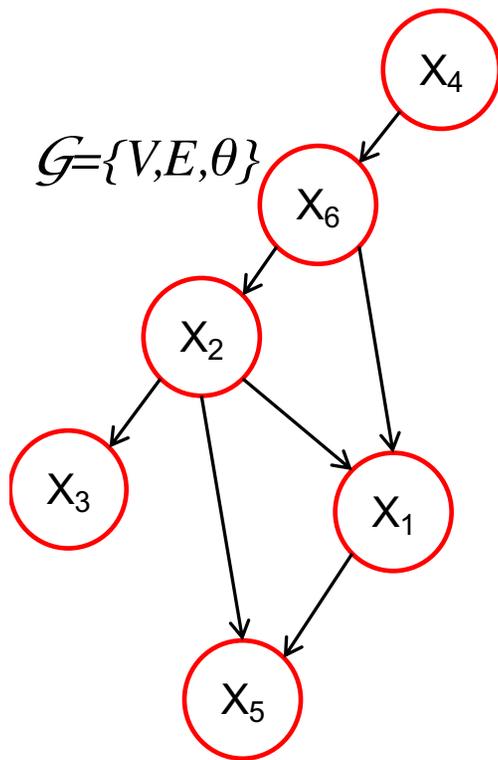
	$X_1 = 0$	$X_1 = 1$	$X_1 = 2$	$X_1 = 3$
$X_2 = 0, X_3 = 0$	0.50	0	0	0.50
$X_2 = 0, X_3 = 1$	0	1.00	0	0
$X_2 = 0, X_3 = 2$	0.18	0.36	0.27	0.18
$X_2 = 0, X_3 = 3$	0.27	0.40	0.30	0.03
$X_2 = 1, X_3 = 0$	0.22	0.45	0.28	0.05
$X_2 = 1, X_3 = 1$	0.43	0	0.28	0.29
$X_2 = 1, X_3 = 2$	NaN	NaN	NaN	NaN
$X_2 = 1, X_3 = 3$	0.39	0.06	0.33	0.22
$X_2 = 2, X_3 = 0$	0.33	0.17	0.33	0.17
$X_2 = 2, X_3 = 1$	0.42	0.11	0.29	0.18
.....

Learning Problems

- Parameters
 - When structure is specified by an expert
 - Experts cannot usually specify parameters
 - Data sets can change over time
 - Need to learn parameters when structure is learnt
- Structure
 - Cannot be easily determined by experts
 - Variables and Structure may change with new data sets

Learning Parameters of BN

- Parameters define local interactions
- Straight-forward since local CPDs



Max Likelihood Estimate

$$P(x_5 | x_1, x_2)$$

	$X_5 = 0$	$X_5 = 1$	$X_5 = 2$	$X_5 = 3$
$X_1 = 0, X_2 = 0$	0.50	0	0	0.50
$X_1 = 0, X_2 = 1$	0	1.00	0	0
$X_1 = 0, X_2 = 2$	0.18	0.36	0.27	0.18
$X_1 = 0, X_2 = 3$	0.27	0.40	0.30	0.03
$X_1 = 0, X_2 = 4$	0.22	0.45	0.28	0.05
$X_1 = 1, X_2 = 0$	0.43	0	0.28	0.29
$X_1 = 1, X_2 = 1$	NaN	NaN	NaN	NaN
$X_1 = 1, X_2 = 2$	0.39	0.06	0.33	0.22
$X_1 = 1, X_2 = 3$	0.33	0.17	0.33	0.17
$X_1 = 1, X_2 = 4$	0.42	0.11	0.29	0.18
.....

Bayesian Estimate
with Dirichlet Prior

	$X_5 = 0$	$X_5 = 1$	$X_5 = 2$	$X_5 = 3$
$X_1 = 0, X_2 = 0$	0.29	0.14	0.29	0.29
$X_1 = 0, X_2 = 1$	0.25	0.25	0.25	0.25
$X_1 = 0, X_2 = 2$	0.25	0.38	0.25	0.12
$X_1 = 0, X_2 = 3$	0.22	0.41	0.31	0.06
$X_1 = 0, X_2 = 4$	0.16	0.52	0.25	0.07
$X_1 = 1, X_2 = 0$	0.29	0.14	0.29	0.29
$X_1 = 1, X_2 = 1$	0.25	0.25	0.25	0.25
$X_1 = 1, X_2 = 2$	0.37	0.05	0.47	0.11
$X_1 = 1, X_2 = 3$	0.33	0.22	0.33	0.11
$X_1 = 1, X_2 = 4$	0.38	0.13	0.29	0.20
.....

Dirichlet Prior

Prior $\theta \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k) \quad \alpha_1 = \dots = \alpha_k = 1$

Likelihood $O = \{o_1, \dots, o_k\} \sim \text{Multinomial}(\theta_1, \dots, \theta_k)$

Posterior $\theta | O \sim \text{Dirichlet}(\alpha'_1, \dots, \alpha'_k)$

$$\alpha'_i = \alpha_i + o_i, \text{ for } i = 1, \dots, k$$

BN Structure Learning

1. Local: Deviance from Independence Tests

$$d_{\chi^2}(\mathcal{D}) = \sum_{x_i, x_j} \frac{(M[x_i, x_j] - M \cdot \hat{P}(x_i) \cdot \hat{P}(x_j))^2}{M \cdot \hat{P}(x_i) \cdot \hat{P}(x_j)} \quad d_I(\mathcal{D}) = \frac{1}{M} \sum_{x_i, x_j} M[x_i, x_j] \log \frac{M[x_i, x_j]}{M[x_i]M[x_j]}$$

1. Rule for accepting/rejecting hypothesis of independence

$$R_{d,t}(\mathcal{D}) = \begin{cases} \text{Accept } d(\mathcal{D}) \leq t \\ \text{Reject } d(\mathcal{D}) > t \end{cases} \quad \begin{array}{l} \text{False Rejection probability due to} \\ \text{choice of } t \text{ is its p-value} \end{array}$$

2. Global: Structure Scoring

- Goodness of Network

Independence Tests

1. For variables x_i, x_j in data set \mathcal{D} of M samples

1. Pearson's Chi-squared (χ^2) statistic

$$d_{\chi^2}(\mathcal{D}) = \sum_{x_i, x_j} \frac{(M[x_i, x_j] - M \cdot \hat{P}(x_i) \cdot \hat{P}(x_j))^2}{M \cdot \hat{P}(x_i) \cdot \hat{P}(x_j)} \quad \text{Sum over all values of } x_i \text{ and } x_j$$

- Independence $\rightarrow d_{\chi^2}(\mathcal{D})=0$, larger value when Joint $M[x,y]$ and expected counts (under independence assumption) differ

2. Mutual Information (K-L divergence) between joint and product of marginals

$$d_I(\mathcal{D}) = \frac{1}{M} \sum_{x_i, x_j} M[x_i, x_j] \log \frac{M[x_i, x_j]}{M[x_i]M[x_j]}$$

- Independence $\rightarrow d_I(\mathcal{D})=0$, otherwise a positive value

• 2. Decision rule

$$R_{d,t}(\mathcal{D}) = \begin{cases} \text{Accept } d(\mathcal{D}) \leq t \\ \text{Reject } d(\mathcal{D}) > t \end{cases}$$

False Rejection probability due to choice of t is its p-value

Structure Scoring

1. Log-likelihood Score for \mathcal{G} with n variables

$$score_L(\mathcal{G} : \mathcal{D}) = \sum_{\mathcal{D}} \sum_{i=1}^n \log \hat{P}(x_i | pa_{x_i}) \quad \text{Sum over all data and variables } x_i$$

2. Bayesian Score

$$score_B(\mathcal{G} : \mathcal{D}) = \log p(\mathcal{D} | \mathcal{G}) + \log p(\mathcal{G})$$

3. Bayes Information Criterion

– With Dirichlet prior over graphs

$$score_{BIC}(\mathcal{G} : \mathcal{D}) = l(\hat{\theta}_G : \mathcal{D}) - \frac{\log M}{2} Dim(\mathcal{G})$$

BN Structure Learning Algorithms

- **Constraint-based**
 - Find structure that best explains determined dependencies
 - Sensitive to errors in testing individual dependencies
 - Koller and Friedman, 2009
- **Score-based**
 - Search the space of networks to find high-scoring structure
 - Since space is super-exponential, need heuristics
 - K2 algorithm((Cooper & Herskovits, 1992)
 - Optimized Branch and Bound (deCampos, Zheng and Ji, 2009)
- **Bayesian Model Averaging**
 - Prediction over all structures
 - May not have closed form, Limitation of X^2
 - Peters, Danzing and Scholkopf, 2011

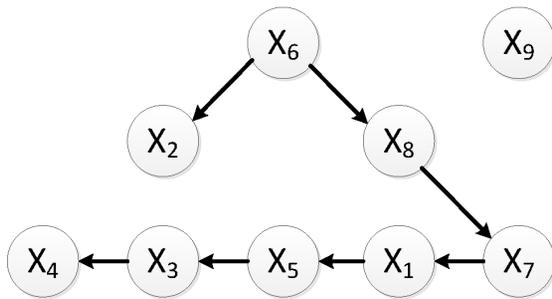
A Causal BN Structure Learning Algorithm*

- Construct PDAG by removing edges from complete undirected graph
- Use χ^2 test to sort dependencies
- Orient most dependent edge using additive noise model
- Apply causal forward propagation to orient other undirected edges
- Repeat until all edges are oriented

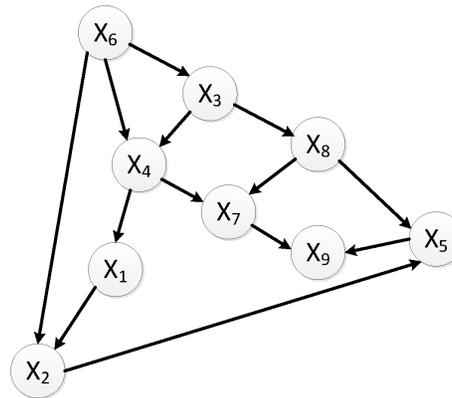
* Zhen and Srihari 2014

Comparison of Algorithms

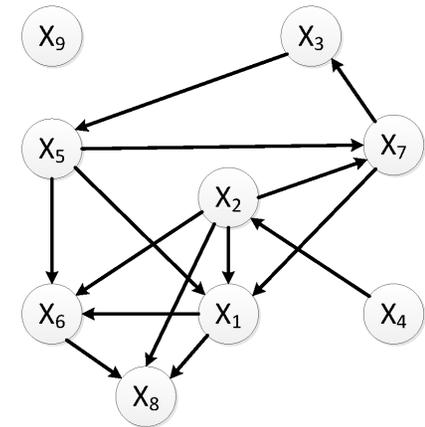
Greedy



B & B



Causal



Data			Algorithm			
Data Set	Type of Data	No. of Vars.	Ind.Vars. (No Edges)	Greedy Algo. [11]	B & B Algo. [10]	Causal Algo.
Set 1	Cursive	9	25994	25329	25642	25228
Set 1	Handprint	9	8059	7898	7301	7094
Set 2	Cursive	12	5316	5142	5139	5008
Set 2	Handprint	12	7825	7004	6976	6956

Some Relevant Papers

- Causality

- <http://www.michaelnielsen.org/ddi/if-correlation-doesnt-imply-causation-then-what-does/>

- Transfer Learning

- Y. Bengio, et al, A Meta-Transfer Objective for Learning to Disentangle Causal Mechanisms, 2019

- Causal Network Learning and Inference

- M. Puri, S. N. Srihari, Y. Tang, "Bayesian Network Structure Learning and Inference Methods for Handwriting," *ICDAR* 2013
- Zhen and Srihari, Learning Causal Networks, 2014
- Srihari, Probabilistic Graphical Models, *Encyclopedia of Social Networks*, Springer 2014
- M. Puri, S. N. Srihari, L. Hanson, "Probabilistic Modeling of Children's Handwriting," *DRR* 2014