

# Structure Learning in Bayesian Networks

Sargur Srihari

[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Topics

- Problem Definition
- Overview of Methods
- Constraint-Based Approaches
- Independence Tests
- Testing for Independence
- Structure Scores

# Problem Assumptions

- We do not know the structure
- Data set is fully observed
  - A strong assumption
- Partially observed data are dealt with using other methods
  - Gradient Ascent and EM
- Assume data  $\mathcal{D}$  is generated i.i.d. from distribution  $\mathcal{P}^*(X)$
- Assume that  $\mathcal{P}^*$  is induced by BN  $\mathcal{G}^*$

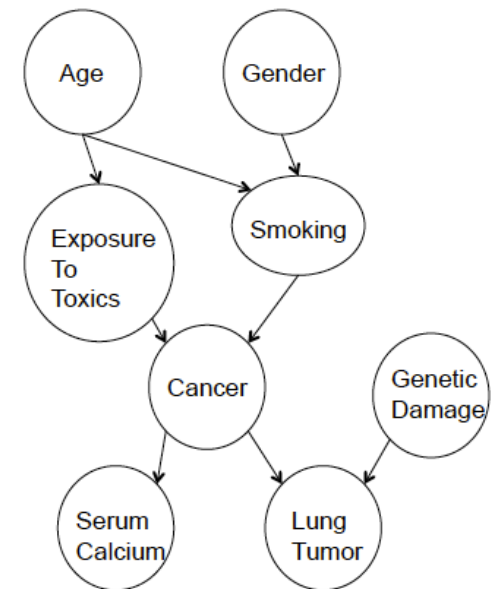
# Need for BN Structure Learning

- Structure

- Causality cannot easily be determined by experts
- Variables and Structure may change with new data sets

- 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



# To what extent do independencies in $G^*$ manifest in $\mathcal{D}$ ?

- Two coins X and Y tossed independently
- We are given data set of 100 instances
- Learn a model for this scenario
- Typical data set:
  - 27 head/head
  - 22 head/tail
  - 25 tail/head
  - 26 tail/tail
- Are the coins independent?



# Coin Tossing Probabilities

- Marginal Probabilities

$$P(X=\text{head})=.49, P(X=\text{tail})=0.51, P(Y=\text{head})=.48, \\ P(Y=\text{tail})=.52$$

- Products of marginals:

- $P(X=\text{head}) \times P(Y=\text{head}) = .49 \times .49 = .24$
- $P(X=\text{head}) \times P(Y=\text{tail}) = .49 \times 0.51 = .25$
- $P(X=\text{tail}) \times P(Y=\text{head}) = .51 \times .48 = .245$
- $P(X=\text{tail}) \times P(Y=\text{tail}) = .51 \times .52 = .26$

- Joint Probabilities

- $P(XY=\text{head-head})=.27$
- $P(XY=\text{head-tail})=.22$
- $P(XY=\text{tail-head})=.25$
- $P(XY=\text{tail-tail})=.26$

- According to empirical distribution: not independent

- But we suspect independence

- since probability of getting exactly 25 in each category is small (approx. 1 in 1,000)



# Rain-Soccer Probabilities

- Scan sports pages for 100 days
- Select an article at random to see
  - If there is mention of rain and soccer
- Marginal Probabilities

$$P(X=\text{rain})=.49, P(X=\text{no rain})=0.51, P(Y=\text{soccer})=.48, P(Y=\text{no soccer})=.52$$

- Joint Probabilities

- $P(XY=\text{rain-soccer})=.27$
- $P(XY=\text{rain-no soccer})=.22$
- $P(XY=\text{no rain-soccer})=.25$
- $P(XY=\text{no rain-no soccer})=.26$

- We suspect there is a weak connection (not independent)
- It is hard to be sure whether the true underlying model has an edge between X and Y



# Goal of Learning

- Correctly constructing network depends on learning goal
- Goal 1: Knowledge Discovery
  - By examining dependencies in the learned network we can learn dependency of variables
    - Can also be done using statistical independence tests
  - Bayesian network reveals much finer structure
    - Distinguish between direct and indirect independencies both of which lead to correlations



# Caution in Knowledge Discovery Goal

- If goal is to understand domain structure best answer we can get is to recover  $\mathcal{G}^*$
- Since there are many I-maps for  $P^*$  we cannot distinguish them from  $\mathcal{D}$
- Thus  $\mathcal{G}^*$  is not *identifiable*
- Best we can do is recover  $\mathcal{G}^*$ 's equivalence class

# Too few or too many edges in $G^*$

- Even learning equivalence class of networks is hard
- Data sampled is noisy
- Need to make decisions about including edges we are less sure about
  - Too few edges means missing out on dependencies
  - Too many edges means spurious dependencies

# Goal of Density Estimation

- Goal is to reason about instances not in our training data
- Network to generalize to new instances
- We want true  $G^*$  which captures true dependencies and independencies
- If we make mistakes, better to include too many edges rather than too few edges
  - With an overly complex structure we can still capture  $\mathcal{P}^*$
- But this also fallacious

# Data Fragmentation with too many edges

- Data set with 20 tosses:

- 3 head/head
- 6 head/tail
- 5 tail/head
- 6 tail/tail



- Marginal Probability

$$P(Y=\text{head})=(3+5)/20 = .4 \quad \text{Based on 20 instances}$$

- Introduce spurious correlation between X and Y

- Conditional Probabilities

$$P(Y=\text{head} \mid X=\text{head})=3/9=1/3 \quad \text{Based on 9 instances}$$

$$P(Y=\text{head} \mid X=\text{tail})=5/11$$

If there was independence  $P(Y/X)=P(Y)$

Data fragmentation causes anomaly

Std dev with  $M$  samples is  $1/\sqrt{M}$  which is .11 for 20 and .17 for 9

# Data Fragmentation with spurious edges

- In a table CPD no of bins grows exponentially with no of parents
- Cost of adding a parent can be very large
- Cost of adding a parent grows with no of parents already there
- It is better to obtain a sparser structure
- We can sometimes learn a better model by learning a model with fewer edges even if it does not represent the true distribution

# Overview of methods

- Three approaches to learning structure

## 1. Constraint-based

- Test for conditional dependence and independence in the data
- Sensitive to failures in independence tests

## 2. Score-based

- BN specifies a statistical model
- Problem is to select best structure
- Hypothesis space is super-exponential  $2^{O(n^2)}$

## 3. Bayesian model averaging

- Generate ensemble of possible structures

# 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

# Elements of BN Structure Learning

1. Local: Independence Tests
  1. Measures of *Deviance*-from-independence between variables
  2. Rule for accepting/rejecting hypothesis of independence
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 ( $\chi^2$ ) 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

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

# 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

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

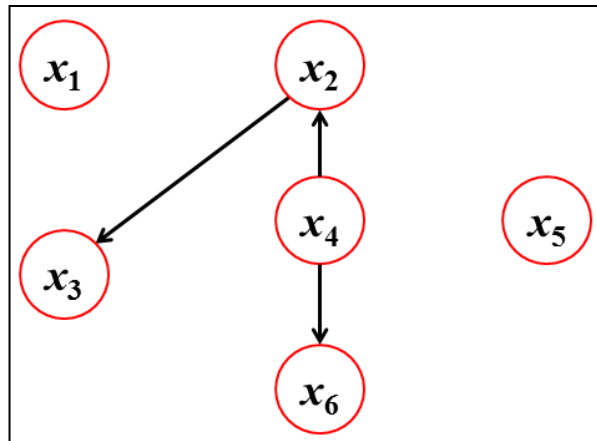
# Heuristic for BN Structure Learning

Consider pairs of variables ordered by  $\chi^2$  value  
Add next edge if score is increased

$$G^* = \{V, E^*, \theta^*\}$$

Start

Score  $s_{k-1}$

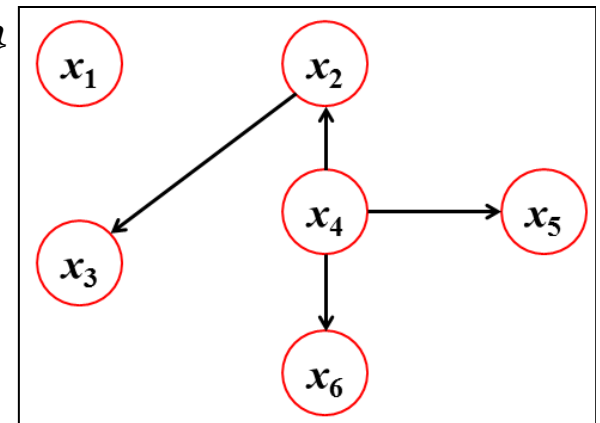


$$G_{c1} = \{V, E_{c1}, \theta_{c1}\}$$

Candidate

$x_4 \rightarrow x_5$

Score  $s_{c1}$

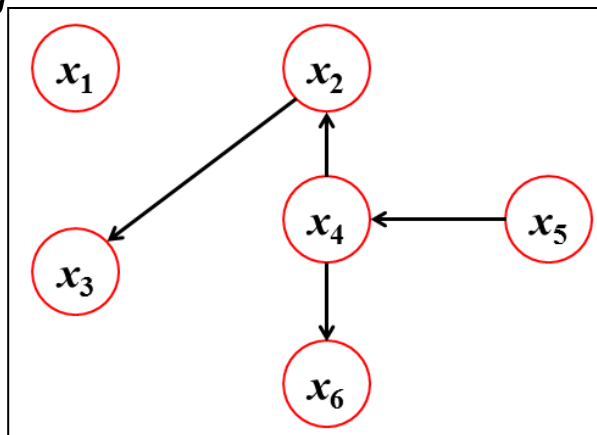


$$G_{c2} = \{V, E_{c2}, \theta_{c2}\}$$

Candidate

$x_5 \rightarrow x_4$

Score  $s_{c2}$



Choose  $G_{c1}$  or  $G_{c2}$   
depending on which one  
increases the score  $s(\mathcal{D}, \mathcal{G})$

evaluate using cross validation  
on validation set