

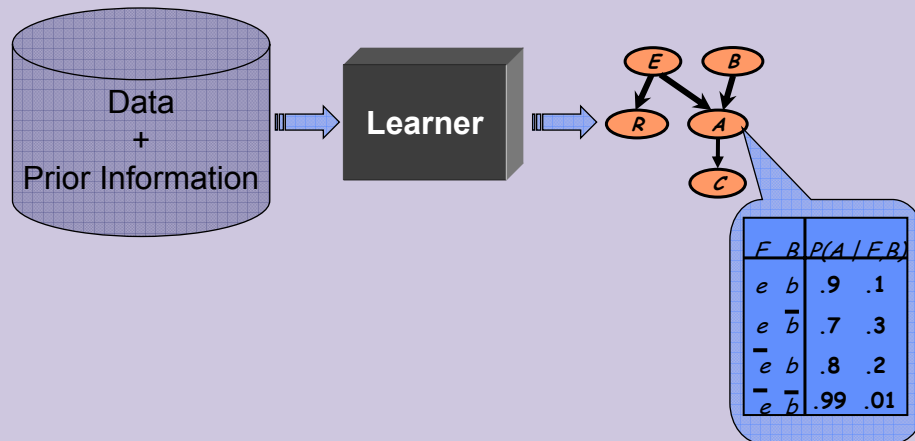
Bayesian Networks Learning

– Parameter Estimation

Learning Bayesian networks

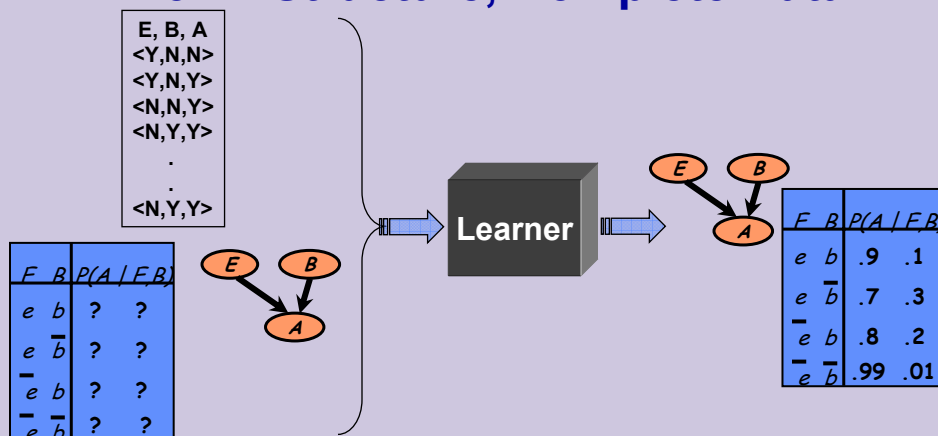
- ◆ Assume the domain $X=\{X_1, \dots, X_n\}$ is governed by some underlying distribution $P^*(X)$
- ◆ P^* is induced by some Bayesian network $B^*=(G, \Theta)$
- ◆ Given a data set $D=\{x[1], \dots, x[M]\}$ of M samples from P^*
- ◆ Samples are *i.i.d. - independent and identically distributed*
- ◆ The task is to recover the Bayesian network model

Learning Bayesian networks



3

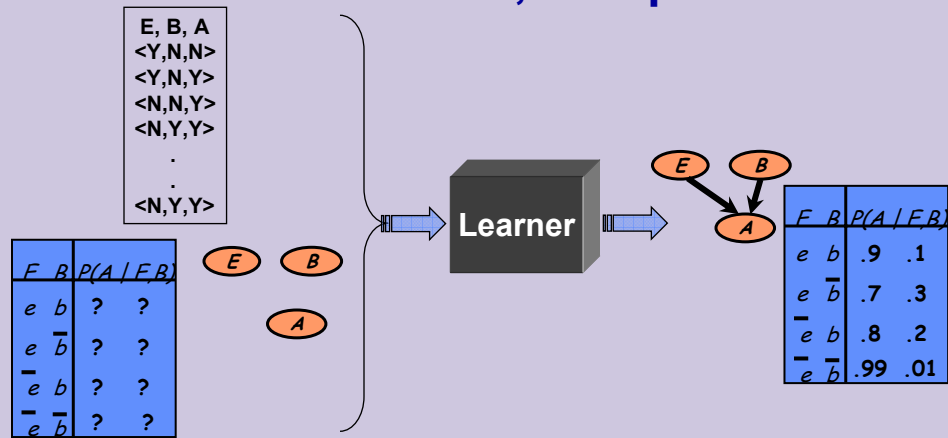
Known Structure, Complete Data



- ◆ Network structure is specified
 - Inducer needs to estimate parameters
- ◆ Data does not contain missing values

4

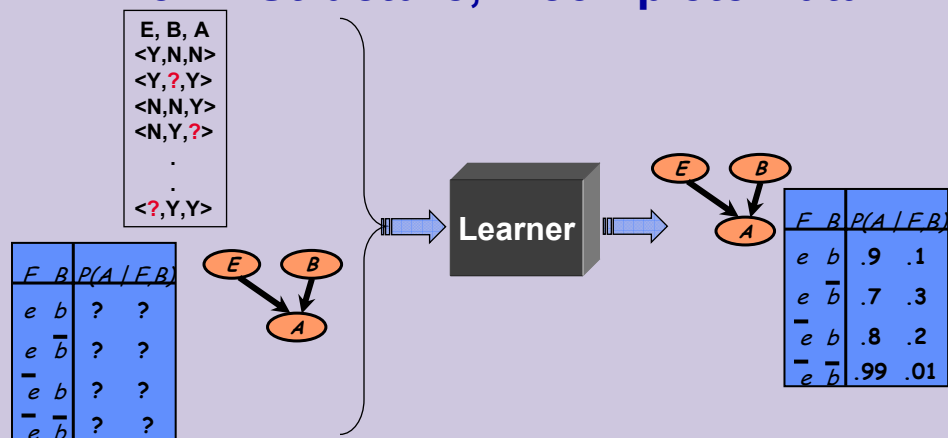
Unknown Structure, Complete Data



- ◆ Network structure is not specified
 - Inducer needs to select arcs & estimate parameters
- ◆ Data does not contain missing values

5

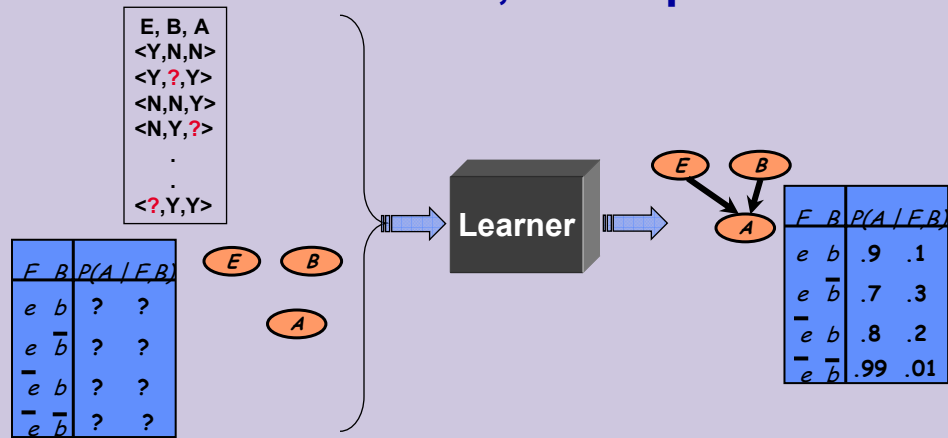
Known Structure, Incomplete Data



- ◆ Network structure is specified
- ◆ Data contains missing values
 - Need to consider assignments to missing values

6

Unknown Structure, Incomplete Data



- ◆ Network structure is not specified
- ◆ Data contains missing values, or hidden variables?
 - Need to consider assignments to missing values

7

Parameter Estimation

- ◆ Network structure is specified
- ◆ Data set D consists of fully observed instances of the network variables

$$D = \{x[1], \dots, x[M]\}$$
- ◆ Estimate network parameters

8

Parameter Estimation

- ◆ Use a set $D=\{x[1], \dots, x[M]\}$ of training samples drawn independently from a parametric model $P(x : \theta)$ to estimate the unknown parameter vector θ
- ◆ Parameter estimation: a classic problem in statistics
 - Maximum-Likelihood (ML) estimation
 - Bayesian estimation

9

Maximum-Likelihood Estimation

- ◆ IID data samples $D=\{x[1], \dots, x[M]\}$
- ◆ *Likelihood function*

$$L(\theta : D) = P(D | \theta) = \prod_{k=1}^M P(x[k] : \theta)$$

the likelihood of θ w.r.t. the set of samples

- ◆ ML estimate of θ is, by definition, the value $\hat{\theta}$ that maximizes $L(\theta:D)$
- ◆ it is usually easier to work with the the *log-likelihood function*

$$l(\theta : D) = \log P(D | \theta) = \sum_{k=1}^M \log P(x[k] : \theta)$$

10

Example: Discrete Case

- ◆ Single binary variable

$$P(X=1) = \theta, \quad P(X=0) = 1 - \theta$$

$$P(x | \theta) = \theta^x (1 - \theta)^{1-x}$$

Bernoulli distribution

$$L(\theta : D) = \prod_{k=1}^M P(x[k] : \theta) = \theta^{N_1} (1 - \theta)^{N_0}$$

- ◆ *Sufficient statistics:*

N_1 : number of 1's in D, N_0 : number of 0's in D

11

- ◆ Log-likelihood

$$l(\theta : D) = N_1 \ln \theta + N_0 \ln (1 - \theta)$$

- ◆ ML estimation

$$\hat{\theta}_{ML} = \frac{N_1}{N_1 + N_0} = \frac{N_1}{M}$$

12

- ◆ Multi-valued discrete random variables $\{1, \dots, K\}$

$$\theta_i = P(X = i)$$

$$P(x | \theta) = \prod_{i=1}^K \theta_i^{\delta_{xi}}$$

$$L(\theta : D) = \prod_{j=1}^M P(x[j] : \theta) = \prod_{i=1}^K \theta_i^{N_i}$$

- Sufficient statistics N_i : the # of times i appears in D

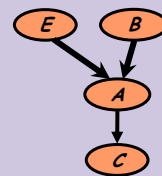
$$\hat{\theta}_{iML} = \frac{N_i}{\sum_j N_j} = \frac{N_i}{M}$$

13

MLE for Bayesian Networks

- ◆ Training data has the form:

$$D = \begin{bmatrix} E[1] & B[1] & A[1] & C[1] \\ \vdots & \vdots & \vdots & \vdots \\ E[M] & B[M] & A[M] & C[M] \end{bmatrix}$$

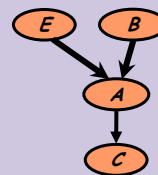


14

Likelihood Function

- ◆ By definition of network, we get

$$\begin{aligned} \mathcal{L}(\Theta : D) &= \prod_m P(E[m], B[m], A[m], C[m] : \Theta) \\ &= \prod_m \left(\begin{array}{l} P(E[m] : \Theta) \\ P(B[m] : \Theta) \\ P(A[m] \mid B[m], E[m] : \Theta) \\ P(C[m] \mid A[m] : \Theta) \end{array} \right) \end{aligned}$$

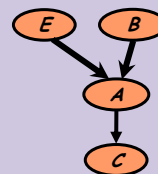


15

Likelihood Function

- ◆ Rewriting terms, we get

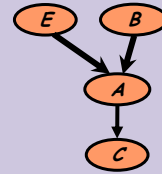
$$\begin{aligned} \mathcal{L}(\Theta : D) &= \prod_m P(E[m], B[m], A[m], C[m] : \Theta) \\ &= \prod_m P(E[m] : \Theta) \\ &\quad \prod_m P(B[m] : \Theta) \\ &\quad \prod_m P(A[m] \mid B[m], E[m] : \Theta) \\ &\quad \prod_m P(C[m] \mid A[m] : \Theta) \end{aligned}$$



16

Likelihood Function

◆ Rewriting terms, we get



$$\begin{aligned}
 L(\Theta : D) &= \prod_m P(E[m], B[m], A[m], C[m] : \Theta) \\
 &= \prod_m P(E[m] : \Theta_E) \\
 &\quad \prod_m P(B[m] : \Theta_B) \\
 &= \prod_m P(A[m] | B[m], E[m] : \Theta_{A|BE}) \\
 &\quad \prod_m P(C[m] | A[m] : \Theta_{C|A})
 \end{aligned}$$

17

General Bayesian Networks

Generalizing for any Bayesian network:

$$\begin{aligned}
 L(\Theta : D) &= \prod_m P(x_1[m], \dots, x_n[m] : \Theta) \\
 &= \prod_i \prod_m P(x_i[m] | Pa_i[m] : \Theta_i) \\
 &= \prod_i L_i(\Theta_i : D) \quad \rightarrow \text{local likelihood function}
 \end{aligned}$$

Global Decomposition of the likelihood function
 \Rightarrow Independent estimation problems

18

MLE

- ◆ Assuming discrete variables (CPTs) leads to further decomposition → *local decomposition of the likelihood function*

$$\begin{aligned}
 L_i(\Theta_i : D) &= \prod_m P(x_i[m] | Pa_i[m] : \Theta_i) \\
 &= \prod_{pa_i} \prod_{m, Pa_i[m]=pa_i} P(x_i[m] | pa_i : \Theta_{x_i|pa_i}) \\
 &= \prod_{pa_i} \prod_{x_i} \theta_{x_i|pa_i}^{N(x_i, pa_i)} \\
 \hat{\theta}_{x_i|pa_i} &= \frac{N(x_i, pa_i)}{N(pa_i)}
 \end{aligned}$$

19

Bayesian Inference

- ◆ Represent uncertainty about parameters using a probability distribution over parameters
- ◆ Learning using Bayes rule

$$P(\theta | x[1], \dots, x[M]) = \frac{P(x[1], \dots, x[M] | \theta) P(\theta)}{P(x[1], \dots, x[M])}$$

Diagram labels:

- Posterior**: $P(\theta | x[1], \dots, x[M])$
- Likelihood**: $P(x[1], \dots, x[M] | \theta)$
- Prior**: $P(\theta)$
- Probability of data**: $P(x[1], \dots, x[M])$

20

Example: Discrete Variable

- ◆ Single binary variable

$$P(X=1) = \theta, \quad P(X=0) = 1 - \theta$$

$$L(\theta : D) = \theta^{N_1} (1 - \theta)^{N_0}$$

- ◆ What prior $p(\theta)$ to use?

21

Beta distribution

$$\text{Beta}(\theta \mid \alpha_1, \alpha_0) = \frac{\Gamma(\alpha_1 + \alpha_0)}{\Gamma(\alpha_1)\Gamma(\alpha_0)} \theta^{\alpha_1-1} (1 - \theta)^{\alpha_0-1}$$

$$0 \leq \theta \leq 1$$

$$E(\theta) = \frac{\alpha_1}{\alpha_1 + \alpha_0}$$

- ◆ The parameters α_1 and α_2 are positive reals, often called *hyperparameters*
- ◆ Gamma Function

$$\Gamma(x+1) = x\Gamma(x)$$

$$\Gamma(1) = 1, \Gamma(x) = (x-1)! \quad \text{for integer } x$$

22

Assume the prior is a Beta distribution

$$p(\theta) = \text{Beta}(\theta | \alpha_1, \alpha_0) = c \theta^{\alpha_1 - 1} (1 - \theta)^{\alpha_0 - 1}$$

The posterior density $p(\theta | D)$

$$\begin{aligned} p(\theta | D) &= c' \cdot p(D | \theta) p(\theta) \\ &= \text{Beta}(\theta | N_1 + \alpha_1, N_0 + \alpha_0) \end{aligned}$$

- ◆ The property that the posterior distribution follows the same parametric form as the prior distribution is called *conjugacy*
- ◆ Beta prior is a *conjugate family* for the binomial distribution

23

$$\begin{aligned} P(X = 1 | D) &= \int P(X = 1 | \theta) p(\theta | D) d\theta \\ &= \int \theta p(\theta | D) d\theta = \frac{N_1 + \alpha_1}{N_1 + N_0 + \alpha_1 + \alpha_0} \equiv \hat{\theta}_{BE} \end{aligned}$$

- ◆ It can be proved that:
If the prior is well-behaved – i.e. does not assign 0 density to any *feasible* parameter value, then both MLE and Bayesian estimate converge to the same value in the limit
- ◆ Both *almost surely* converge to the underlying distribution $P(X)$
- ◆ But the ML and Bayesian approaches behave differently when the number of samples is small

24

- ◆ Multi-valued discrete random variables $\{1, \dots, k\}$

$$\theta_i = P(X = i)$$

$$P(D | \theta) = \prod_{i=1}^k \theta_i^{N_i}$$

Sufficient statistics N_i : the # of times i appears in D

- ◆ Assume the prior $p(\theta)$ is a Dirichlet distribution $\text{Dir}(\theta | \alpha)$

25

Dirichlet distribution with hyperparameters α_i 's

$$\text{Dir}(\theta | \alpha) = \frac{\Gamma(\alpha)}{\prod_{i=1}^k \Gamma(\alpha_i)} \prod_{i=1}^k \theta_i^{\alpha_i - 1}$$

$$0 \leq \theta_i \leq 1, \quad \sum_{i=1}^k \theta_i = 1, \quad \alpha = \sum_{i=1}^k \alpha_i$$

$$E(\theta_i) = \frac{\alpha_i}{\alpha}$$

26

- ◆ Multi-valued discrete random variables $\{1, \dots, k\}$
- ◆ Assume the prior $p(\theta)$ is a Dirichlet distribution $\text{Dir}(\theta|\alpha)$ with *hyperparameters* α_i 's
- ◆ Then the posterior density $p(\theta | D)$ is also a Dirichlet distribution with hyperparameters $\alpha_1 + N_1, \dots, \alpha_k + N_k$

$$\begin{aligned} p(\boldsymbol{\theta} | D) &= c \cdot P(D | \boldsymbol{\theta}) p(\boldsymbol{\theta}) \\ &= \text{Dir}(\boldsymbol{\theta} | \mathbf{N} + \boldsymbol{\alpha}) \end{aligned}$$

- ◆ Dirichlet prior is a conjugate family for the multinomial distribution

27

- ◆ Bayesian estimates

$$\begin{aligned} P(X = i | D) &= \int P(X = i | \boldsymbol{\theta}) p(\boldsymbol{\theta} | D) d\boldsymbol{\theta} \\ &= \int \theta_i \text{Dir}(\boldsymbol{\theta} | \mathbf{N} + \boldsymbol{\alpha}) d\boldsymbol{\theta} = \frac{N_i + \alpha_i}{M + \alpha} \equiv \hat{\theta}_{iBE} \end{aligned}$$

- ◆ The hyperparameters α_i can be thought of as “*imaginary*” counts from our prior experience
- ◆ α : imaginary *equivalent sample size*
- ◆ Let p_i be prior belief about θ_i : $\alpha_i = \alpha p_i$
- ◆ The larger the equivalent sample size, the more confident we are in our prior
- ◆ *Laplace estimates*: $\alpha = k$, $\alpha_i = 1$

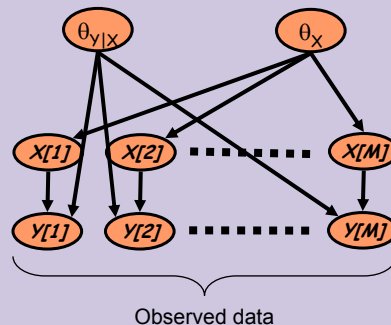
28

Summary of Bayesian estimation

- ◆ Treat the unknown parameters as random variables
- ◆ Assume a prior distribution for the unknown parameters
- ◆ Update the distribution of the parameters based on data
- ◆ Finally compute $p(x|D)$

29

Bayesian Estimation in BNs



- ◆ **Meta-network** for $P(\Theta, D)$
- ◆ Priors for each parameter group are independent
- ◆ Data instances are independent given the unknown parameters

30

Bayesian Estimation in BNs

- ◆ Global parameter independence assumption

$$P(\Theta) = \prod_i P(\Theta_{X_i | Pa_i})$$

- ◆ Global Decomposition of the likelihood function

$$L(\Theta : D) = P(D | \Theta) = \prod_i L_i(\Theta_{X_i | Pa_i} : D)$$

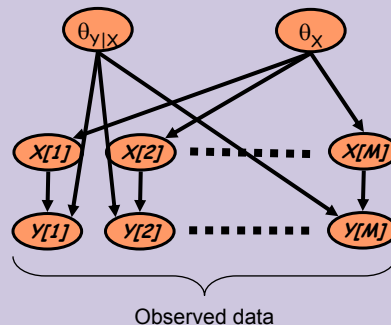
- ◆ Posterior parameter independence

$$P(\Theta | D) = \prod_i P(\Theta_{X_i | Pa_i} | D)$$

- ◆ We can solve the prediction problem for each CPD independently

31

Bayesian Estimation in BNs



- ◆ This can be “read” from the network by d-separation

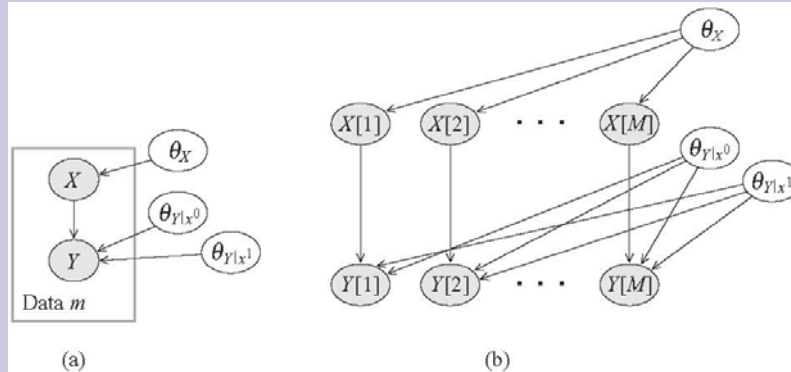
Complete data \Rightarrow

posteriors on parameters are independent

- ◆ Can compute posterior over parameters separately!

32

Bayesian Estimation in BNs - CPTs



- ◆ Local parameter independence assumption

$$P(\Theta_{X_i|pa_i}) = \prod_{pa_i} P(\Theta_{X_i|pa_i})$$

33

Bayesian Nets & Bayesian Prediction

- ◆ Posterior parameter independence

$$P(\Theta | D) = \prod_i \prod_{pa_i} P(\Theta_{X_i|pa_i} | D)$$

- ◆ Assume Dirichlet prior

$$P(\Theta_{X_i|pa_i}) = Dir(\Theta_{X_i|pa_i} | \alpha_{x_i|pa_i}, \dots)$$

- ◆ Then

$$P(\Theta_{X_i|pa_i} | D) = Dir(\Theta_{X_i|pa_i} | \alpha_{x_i|pa_i} + N(x_i, pa_i), \dots)$$

34

Bayesian Nets & Bayesian Prediction

◆ Bayesian estimation

$$\begin{aligned}\tilde{\theta}_{x_i|pa_i} &= P(X_i = x_i \mid Pa_i = pa_i, D) \\ &= \frac{\alpha_{x_i|pa_i} + N(x_i, pa_i)}{\alpha_{pa_i} + N(pa_i)}\end{aligned}$$

35

Assessing Priors for Bayesian Nets

The BDe prior

- ◆ Introduce an equivalent sample size α and a prior distribution $P'(X)$, and set

$$\alpha_{x_i|pa_i} = \alpha P'(x_i, pa_i)$$

- ◆ We can represent P' as a BN (G_0, Θ_0) , and set

$$\alpha_{x_i|pa_i} = \alpha P(x_i, pa_i^G \mid G_0, \Theta_0)$$

Use BN inference to compute this

- ◆ E.g., empty BN with uniform distribution

$$\alpha_{x_i|pa_i} = \frac{\alpha}{|X_i| \prod |Pa_i|}$$

36

Learning Parameters: Summary

- ◆ Estimation relies on **sufficient statistics**

- For multinomials: counts $N(x_i, pa_i)$
- Parameter estimation

$$\hat{\theta}_{x_i|pa_i} = \frac{N(x_i, pa_i)}{N(pa_i)} \quad \tilde{\theta}_{x_i|pa_i} = \frac{\alpha_{x_i|pa_i} + N(x_i, pa_i)}{\alpha_{pa_i} + N(pa_i)}$$

MLE Bayesian (Dirichlet)

- ◆ Both are asymptotically equivalent and consistent
- ◆ Both can be implemented in an on-line manner by accumulating sufficient statistics