

Learning Bayesian Networks from Data

Nir Friedman
Hebrew U.

Daphne Koller
Stanford

Overview

- ◆ Introduction
- ◆ Parameter Estimation
- ◆ Model Selection
- ◆ Structure Discovery
- ◆ Incomplete Data
- ◆ Learning from Structured Data

2

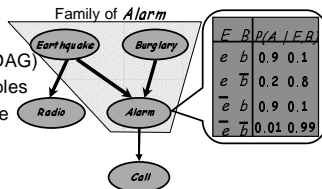
Bayesian Networks

Compact representation of probability distributions via conditional independence

Qualitative part:

Directed acyclic graph (DAG)

- ◆ Nodes - random variables
- ◆ Edges - direct influence



Together:

Define a unique distribution in a factored form

Quantitative part:
Set of conditional probability distributions

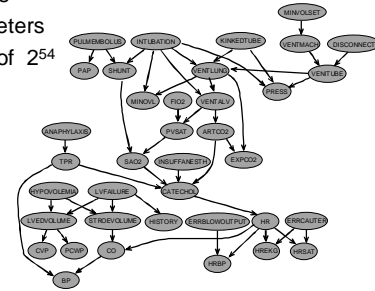
$$P(B, E, A, C, R) = P(B)P(E)P(A | B, E)P(R | E)P(C | A)$$

3

Example: "ICU Alarm" network

Domain: Monitoring Intensive-Care Patients

- ◆ 37 variables
- ◆ 509 parameters
- ... instead of 2^{54}



4

Inference

◆ Posterior probabilities

- Probability of any event given any evidence

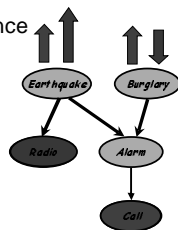
◆ Most likely explanation

- Scenario that explains evidence

◆ Rational decision making

- Maximize expected utility
- Value of Information

◆ Effect of intervention



5

Why learning?

Knowledge acquisition bottleneck

- ◆ Knowledge acquisition is an expensive process
- ◆ Often we don't have an expert

Data is cheap

- ◆ Amount of available information growing rapidly
- ◆ Learning allows us to construct models from raw data

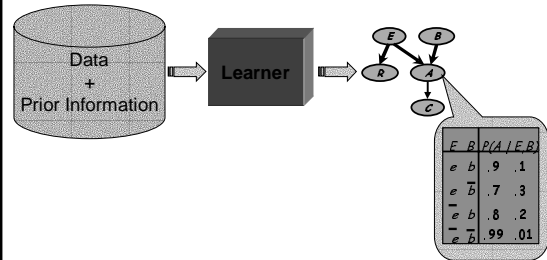
6

Why Learn Bayesian Networks?

- ◆ Conditional independencies & graphical language capture structure of many real-world distributions
- ◆ Graph structure provides much insight into domain
 - Allows "knowledge discovery"
- ◆ Learned model can be used for many tasks
- ◆ Supports all the features of probabilistic learning
 - Model selection criteria
 - Dealing with missing data & hidden variables

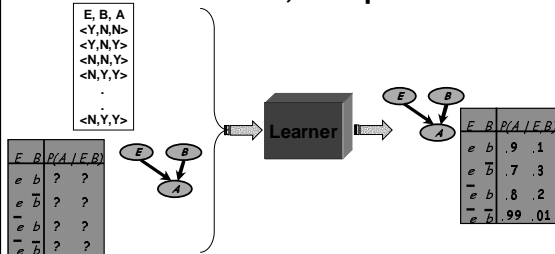
7

Learning Bayesian networks



8

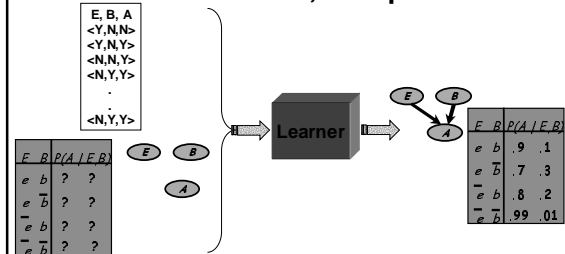
Known Structure, Complete Data



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

9

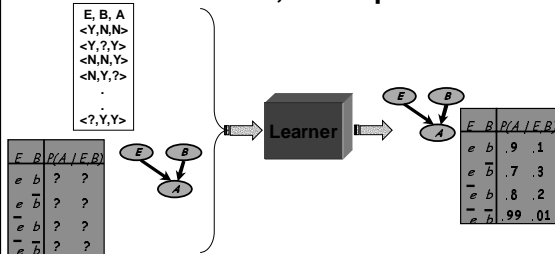
Unknown Structure, Complete Data



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

10

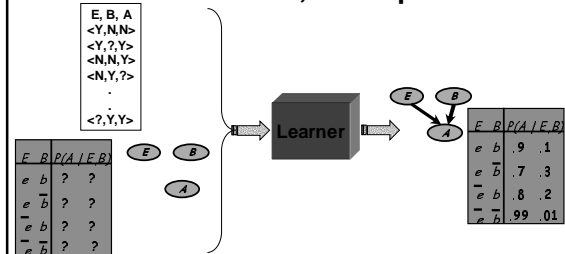
Known Structure, Incomplete Data



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

11

Unknown Structure, Incomplete Data



- ◆ Network structure is not specified
- ◆ Data contains missing values
 - Need to consider assignments to missing values

12

Overview

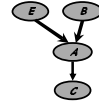
- ◆ Introduction
- ◆ **Parameter Estimation**
 - Likelihood function
 - Bayesian estimation
- ◆ Model Selection
- ◆ Structure Discovery
- ◆ Incomplete Data
- ◆ Learning from Structured Data

13

Learning Parameters

- ◆ Training data has the form:

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

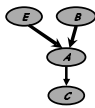


14

Likelihood Function

- ◆ Assume i.i.d. samples
- ◆ Likelihood function is

$$L(\theta : D) = \prod_m P(E[m], B[m], A[m], C[m] : \theta)$$



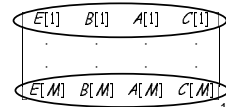
15

Likelihood Function

- ◆ By definition of network, we get

$$L(\theta : D) = \prod_m P(E[m], B[m], A[m], C[m] : \theta)$$

$$= \prod_m \begin{pmatrix} P(E[m] : \theta) \\ P(B[m] : \theta) \\ P(A[m] | B[m], E[m] : \theta) \\ P(C[m] | A[m] : \theta) \end{pmatrix}$$



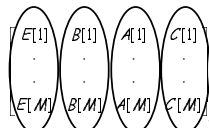
16

Likelihood Function

- ◆ Rewriting terms, we get

$$L(\theta : D) = \prod_m P(E[m], B[m], A[m], C[m] : \theta)$$

$$= \prod_m P(E[m] : \theta) \prod_m P(B[m] : \theta) \prod_m P(A[m] | B[m], E[m] : \theta) \prod_m P(C[m] | A[m] : \theta)$$



17

General Bayesian Networks

Generalizing for any Bayesian network:

$$L(\theta : D) = \prod_m P(x_1[m], \dots, x_n[m] : \theta)$$

$$= \prod_i \prod_m P(x_i[m] | Pa_i[m] : \theta,)$$

$$= \prod_i L_i(\theta, : D)$$

Decomposition

⇒ Independent estimation problems

18

Likelihood Function: Multinomials

$L(\theta : D) = P(D | \theta) = \prod_m P(x[m] | \theta)$

- The likelihood for the sequence H, T, T, H, H is

$L(\theta : D) = \theta \cdot (1 - \theta) \cdot (1 - \theta) \cdot \theta \cdot \theta$

General case: $L(\theta : D) = \prod_{k=1}^K \theta_k^{N_k}$

- Count of k^{th} outcome in D
- Probability of k^{th} outcome

19

Bayesian Inference

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

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

Labels: Likelihood, Prior, Posterior, Probability of data

20

Bayesian Inference

- Represent Bayesian distribution as Bayes net

- The values of X are independent given θ
- $P(x[m] | \theta) = \theta$
- Bayesian prediction is inference in this network

21

Example: Binomial Data

- Prior: uniform for θ in $[0, 1]$
- $\Rightarrow P(\theta | D) \propto$ the likelihood $L(\theta : D)$

$$P(\theta | x[1], \dots, x[M]) \propto P(x[1], \dots, x[M] | \theta) \cdot P(\theta)$$

$(N_H, N_T) = (4, 1)$

- MLE for $P(X = H)$ is $4/5 = 0.8$
- Bayesian prediction is

$$P(x[M+1] = H | D) = \int \theta \cdot P(\theta | D) d\theta = \frac{5}{7} = 0.7142\dots$$

22

Dirichlet Priors

- Recall that the likelihood function is

$$L(\theta : D) = \prod_{k=1}^K \theta_k^{N_k}$$

- Dirichlet prior with hyperparameters $\alpha_1, \dots, \alpha_K$

$$P(\theta) \propto \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

\Rightarrow the posterior has the same form, with hyperparameters $\alpha_1 + N_1, \dots, \alpha_K + N_K$

$$P(\theta | D) \propto P(\theta) P(D | \theta) \propto \prod_{k=1}^K \theta_k^{\alpha_k - 1} \prod_{k=1}^K \theta_k^{N_k} = \prod_{k=1}^K \theta_k^{\alpha_k + N_k - 1}$$

23

Dirichlet Priors - Example

24

Dirichlet Priors (cont.)

◆ If $P(\theta)$ is Dirichlet with hyperparameters $\alpha_1, \dots, \alpha_K$

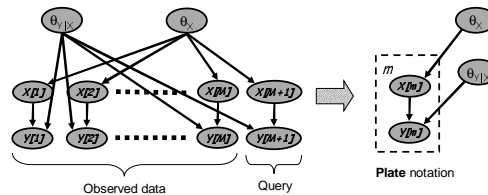
$$P(X[1] = k) = \int \theta_k \cdot P(\theta) d\theta = \frac{\alpha_k}{\sum_{\ell} \alpha_{\ell}}$$

◆ Since the posterior is also Dirichlet, we get

$$P(X[M+1] = k | D) = \int \theta_k \cdot P(\theta | D) d\theta = \frac{\alpha_k + N_k}{\sum_{\ell} (\alpha_{\ell} + N_{\ell})}$$

25

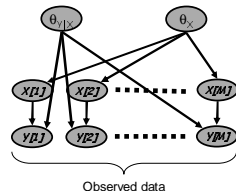
Bayesian Nets & Bayesian Prediction



- ◆ Priors for each parameter group are independent
- ◆ Data instances are independent given the unknown parameters

26

Bayesian Nets & Bayesian Prediction



◆ We can also "read" from the network:

Complete data \Rightarrow

posteriors on parameters are independent

◆ Can compute posterior over parameters separately!

27

Learning Parameters: Summary

◆ Estimation relies on **sufficient statistics**

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

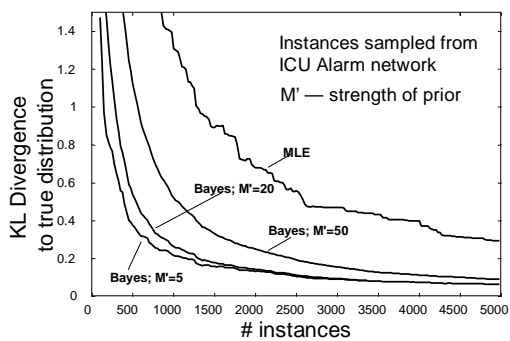
$$\hat{\theta}_{x_i, p a_i} = \frac{N(x_i, p a_i)}{N(p a_i)} \quad \tilde{\theta}_{x_i, p a_i} = \frac{\alpha(x_i, p a_i) + N(x_i, p a_i)}{\alpha(p a_i) + N(p a_i)}$$

MLE Bayesian (Dirichlet)

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

28

Learning Parameters: Case Study



29

Overview

- ◆ Introduction
- ◆ Parameter Learning
- ◆ **Model Selection**
 - Scoring function
 - Structure search
- ◆ Structure Discovery
- ◆ Incomplete Data
- ◆ Learning from Structured Data

30

Why Struggle for Accurate Structure?

Missing an arc

- ◆ Cannot be compensated for by fitting parameters
- ◆ Wrong assumptions about domain structure

Adding an arc

- ◆ Increases the number of parameters to be estimated
- ◆ Wrong assumptions about domain structure

31

Score-based Learning

Define scoring function that evaluates how well a structure matches the data

Search for a structure that maximizes the score

32

Likelihood Score for Structure

$$\ell(\mathcal{G} : D) = \log L(\mathcal{G} : D) = M \sum_i (I(X_i; Pa_i^{\mathcal{G}}) - H(X_i))$$

Mutual information between X_i and its parents

- ◆ Larger dependence of X_i on $Pa_i \Rightarrow$ higher score
- ◆ Adding arcs always helps
 - $I(X; Y) \leq I(X; \{Y, Z\})$
 - Max score attained by fully connected network
 - Overfitting: A bad idea...

33

Bayesian Score

Likelihood score: $L(\mathcal{G} : D) = P(D | \mathcal{G}, \hat{\theta}_{\mathcal{G}})$

Bayesian approach:

- ◆ Deal with uncertainty by assigning probability to all possibilities

$$P(D | \mathcal{G}) = \int P(D | \mathcal{G}, \theta) P(\theta | \mathcal{G}) d\theta$$

Marginal Likelihood Likelihood Prior over parameters

$$P(\mathcal{G} | D) = \frac{P(D | \mathcal{G}) P(\mathcal{G})}{P(D)}$$

34

Marginal Likelihood: Multinomials

Fortunately, in many cases integral has closed form

- ◆ $P(\theta)$ is Dirichlet with hyperparameters $\alpha_1, \dots, \alpha_K$
- ◆ D is a dataset with sufficient statistics N_1, \dots, N_K

Then

$$P(D) = \frac{\Gamma(\sum_{\ell} \alpha_{\ell})}{\Gamma(\sum_{\ell} (\alpha_{\ell} + N_{\ell}))} \prod_{\ell} \frac{\Gamma(\alpha_{\ell} + N_{\ell})}{\Gamma(\alpha_{\ell})}$$

35

Marginal Likelihood: Bayesian Networks

- ◆ Network structure determines form of marginal likelihood

		1	2	3	4	5	6	7
x	H	T	T	H	T	H	H	
y	H	T	H	H	T	T	H	

Network 1:
Two Dirichlet marginal likelihoods

$P(\theta_x)$

) Integral over θ_x

$P(\theta_y)$

) Integral over θ_y

36

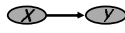
Marginal Likelihood: Bayesian Networks

◆ Network structure determines form of marginal likelihood

	1	2	3	4	5	6	7
x	H	T	T	H	T	H	H
y	H	T	H	H	T	T	H

Network 2:

Three Dirichlet marginal likelihoods



$P(\dots)$) Integral over θ_X

$P(\dots)$) Integral over $\theta_{Y|X=H}$

$P(\dots)$) Integral over $\theta_{Y|X=T}$

37

Marginal Likelihood for Networks

The marginal likelihood has the form:

$$P(D | \mathcal{G}) = \prod_i \prod_{pa^c} \text{Dirichlet marginal likelihood for multinomial } P(X_i | pa_i)$$

$$\frac{\Gamma(\alpha(pa^c))}{\Gamma(\alpha(pa^c) + N(pa^c))} \prod_x \frac{\Gamma(\alpha(x, pa^c) + N(x, pa^c))}{\Gamma(\alpha(x, pa^c))}$$

$N(\dots)$ are counts from the data
 $\alpha(\dots)$ are hyperparameters for each family given \mathcal{G}

38

Bayesian Score: Asymptotic Behavior

$$\log P(D | \mathcal{G}) = \ell(\mathcal{G} : D) - \frac{\log M}{2} \dim(\mathcal{G}) + O(1)$$

$$= M \sum_i (\underbrace{I(X_i; Pa_i^c)}_{\text{Fit dependencies in empirical distribution}}) - \underbrace{\frac{\log M}{2} \dim(\mathcal{G})}_{\text{Complexity penalty}} + O(1)$$

Fit dependencies in empirical distribution

Complexity penalty

- ◆ As M (amount of data) grows,
 - Increasing pressure to fit dependencies in distribution
 - Complexity term avoids fitting noise
- ◆ Asymptotic equivalence to MDL score
- ◆ Bayesian score is **consistent**
 - Observed data eventually overrides prior

39

Structure Search as Optimization

Input:

- Training data
- Scoring function
- Set of possible structures

Output:

- A network that maximizes the score

Key Computational Property: Decomposability:

$$\text{score}(\mathcal{G}) = \sum \text{score}(\text{family of } X \text{ in } \mathcal{G})$$

40

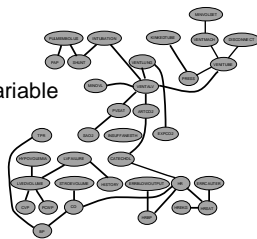
Tree-Structured Networks

Trees:

- ◆ At most one parent per variable

Why trees?

- ◆ Elegant math
 - = we can solve the optimization problem
- ◆ Sparse parameterization
 - = avoid overfitting



41

Learning Trees

- ◆ Let $p(i)$ denote parent of X_i
- ◆ We can write the Bayesian score as

$$\text{Score}(\mathcal{G} : D) = \sum_i \text{Score}(X_i : Pa_i)$$

$$= \sum_i (\text{Score}(X_i : X_{p(i)}) - \text{Score}(X_i)) + \sum_i \text{Score}(X_i)$$

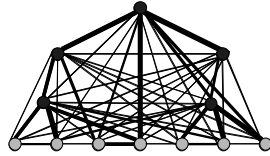
Improvement over "empty" network

Score of "empty" network

Score = sum of edge scores + constant

42

Learning Trees



- ◆ Set $w(j \rightarrow i) = \text{Score}(X_j \rightarrow X_i) - \text{Score}(X_i)$
- ◆ Find tree (or forest) with maximal weight
 - Standard max spanning tree algorithm — $O(n^2 \log n)$

Theorem: This procedure finds tree with max score

43

Beyond Trees

When we consider more complex network, the problem is not as easy

- ◆ Suppose we allow at most two parents per node
- ◆ A greedy algorithm is no longer guaranteed to find the optimal network

- ◆ In fact, no efficient algorithm exists

Theorem: Finding maximal scoring structure with at most k parents per node is NP-hard for $k > 1$

44

Heuristic Search

- ◆ Define a search space:
 - search states are possible structures
 - operators make small changes to structure
- ◆ Traverse space looking for high-scoring structures
- ◆ Search techniques:
 - Greedy hill-climbing
 - Best first search
 - Simulated Annealing
 - ...

45

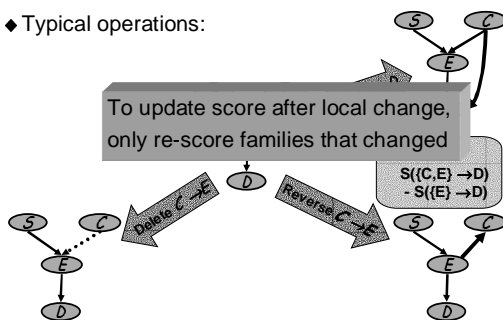
Local Search

- ◆ Start with a given network
 - empty network
 - best tree
 - a random network
- ◆ At each iteration
 - Evaluate all possible changes
 - Apply change based on score
- ◆ Stop when no modification improves score

46

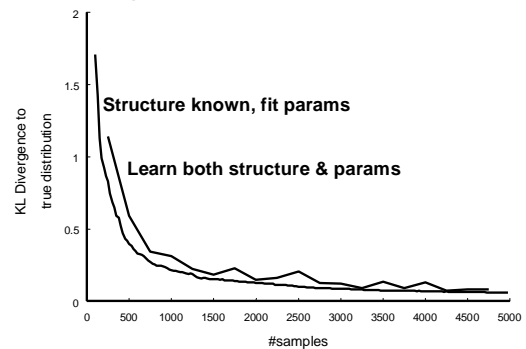
Heuristic Search

- ◆ Typical operations:



47

Learning in Practice: Alarm domain



48

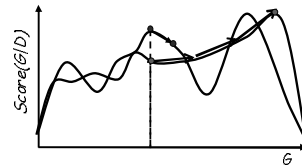
Local Search: Possible Pitfalls

- ◆ Local search can get stuck in:
 - **Local Maxima:**
 - All one-edge changes reduce the score
 - **Plateaux:**
 - Some one-edge changes leave the score unchanged
- ◆ Standard heuristics can escape both
 - Random restarts
 - TABU search
 - Simulated annealing

49

Improved Search: Weight Annealing

- ◆ Standard annealing process:
 - Take bad steps with probability $\propto \exp(\Delta \text{score}/t)$
 - Probability increases with temperature
- ◆ Weight annealing:
 - Take uphill steps relative to **perturbed** score
 - Perturbation increases with temperature



50

Perturbing the Score

- ◆ Perturb the score by reweighting instances
- ◆ Each weight sampled from distribution:
 - Mean = 1
 - Variance \propto temperature
- ◆ Instances sampled from "original" distribution
- ◆ ... but perturbation changes emphasis

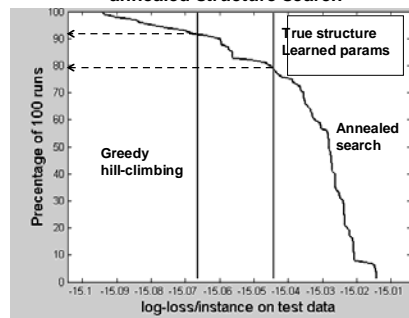
Benefit:

- ◆ **Allows global moves in the search space**

51

Weight Annealing: ICU Alarm network

Cumulative performance of 100 runs of annealed structure search



52

Structure Search: Summary

- ◆ Discrete optimization problem
- ◆ In some cases, optimization problem is easy
 - Example: learning trees
- ◆ In general, NP-Hard
 - Need to resort to heuristic search
 - In practice, search is relatively fast (~100 vars in ~2-5 min):
 - Decomposability
 - Sufficient statistics
 - Adding randomness to search is critical

53

Overview

- ◆ Introduction
- ◆ Parameter Estimation
- ◆ Model Selection
- ◆ **Structure Discovery**
- ◆ Incomplete Data
- ◆ Learning from Structured Data

54

Structure Discovery

Task: Discover structural properties

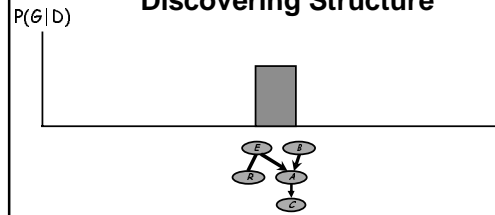
- Is there a direct connection between X & Y
- Does X separate between two “subsystems”
- Does X causally effect Y

Example: scientific data mining

- Disease properties and symptoms
- Interactions between the expression of genes

55

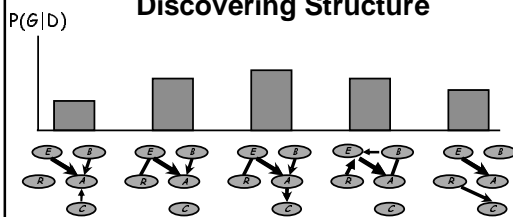
Discovering Structure



- ◆ Current practice: model selection
 - Pick a single high-scoring model
 - Use that model to infer domain structure

56

Discovering Structure



Problem

- Small sample size \Rightarrow many high scoring models
- Answer based on one model often useless
- Want features common to many models

57

Bayesian Approach

- ◆ Posterior distribution over structures
- ◆ Estimate probability of **features**
 - Edge $X \rightarrow Y$
 - Path $X \rightarrow \dots \rightarrow Y$
 - ...

$$P(f | D) = \sum_G f(G) P(G | D)$$

Bayesian score for G

Feature of G_i , e.g., $X \rightarrow Y$

Indicator function for feature f

58

MCMC over Networks

- ◆ Cannot enumerate structures, so sample structures

$$P(f(G) | D) \approx \frac{1}{n} \sum_{i=1}^n f(G_i)$$

- ◆ MCMC Sampling

- Define Markov chain over BNs
- Run chain to get samples from posterior $P(G | D)$

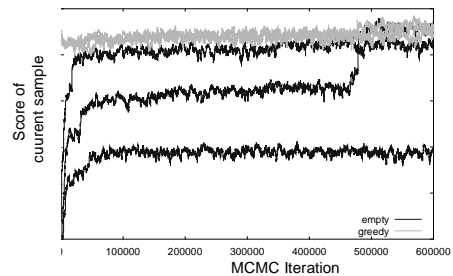
Possible pitfalls:

- Huge (superexponential) number of networks
- Time for chain to converge to posterior is unknown
- Islands of high posterior, connected by low bridges

59

ICU Alarm BN: No Mixing

- ◆ 500 instances:



- ◆ The runs clearly do not mix

60

Effects of Non-Mixing

- ◆ Two MCMC runs over same 500 instances
- ◆ Probability estimates for edges for two runs

Probability estimates highly variable, nonrobust

61

Fixed Ordering

Suppose that

- ◆ We know the **ordering** of variables
 - say, $X_1 \succ X_2 \succ X_3 \succ X_4 \succ \dots \succ X_n$
- ◆ Limit number of parents per nodes to k

$2^{k \cdot n \cdot \log n}$ networks

Intuition: Order **decouples** choice of parents

- ◆ Choice of $Pa(X_i)$ does not restrict choice of $Pa(X_{i+1})$

Upshot: Can compute efficiently in closed form

- ◆ Likelihood $P(D | \prec)$
- ◆ Feature probability $P(f | D, \prec)$

62

Our Approach: Sample Orderings

We can write

$$P(f | D) = \sum_{\prec} P(f | \prec, D) P(\prec | D)$$

Sample orderings and approximate

$$P(f | D) \approx \sum_{j=1}^n P(f | \prec_j, D)$$

- ◆ MCMC Sampling
 - Define Markov chain over orderings
 - Run chain to get samples from posterior $P(\prec | D)$

63

Mixing with MCMC-Orderings

- ◆ 4 runs on ICU-Alarm with 500 instances
 - fewer iterations than MCMC-Nets
 - approximately same amount of computation

Process appears to be mixing!

64

Mixing of MCMC runs

- ◆ Two MCMC runs over same instances
- ◆ Probability estimates for edges

Probability estimates very robust

65

Application: Gene expression

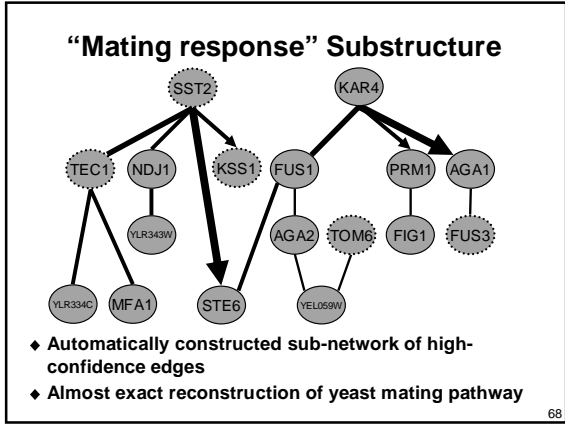
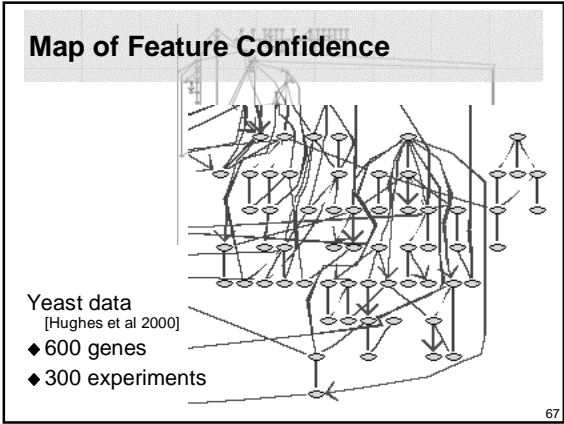
Input:

- ◆ Measurement of gene expression under different conditions
 - Thousands of genes
 - Hundreds of experiments

Output:

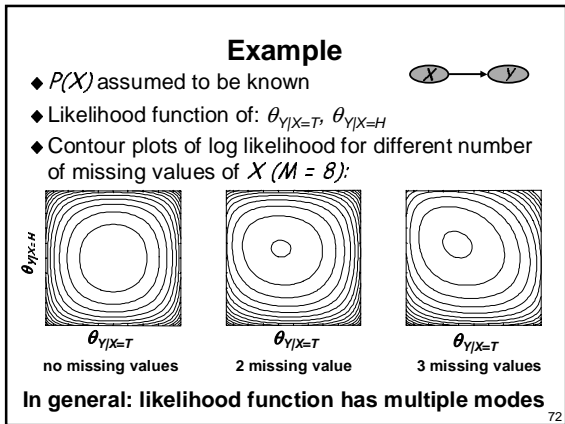
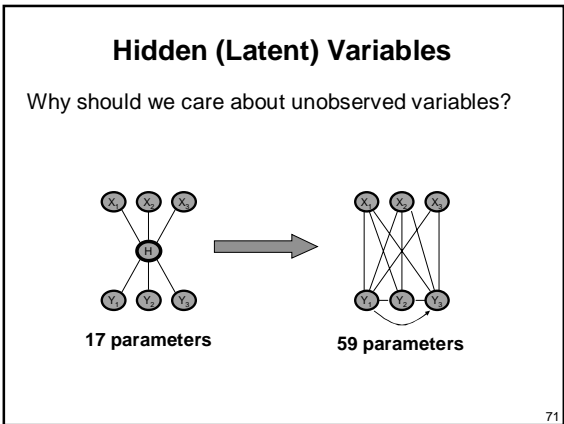
- ◆ Models of gene interaction
 - Uncover pathways

66



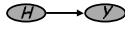
- ### Overview
- ◆ Introduction
 - ◆ Parameter Estimation
 - ◆ Model Selection
 - ◆ Structure Discovery
 - ◆ **Incomplete Data**
 - Parameter estimation
 - Structure search
 - ◆ Learning from Structured Data
- 69

- ### Incomplete Data
- Data is often **incomplete**
- ◆ Some variables of interest are not assigned values
- This phenomenon happens when we have
- ◆ **Missing values:**
 - Some variables unobserved in some instances
 - ◆ **Hidden variables:**
 - Some variables are never observed
 - We might not even know they exist
- 70



Incomplete Data

- ◆ In the presence of incomplete data, the likelihood can have multiple maxima

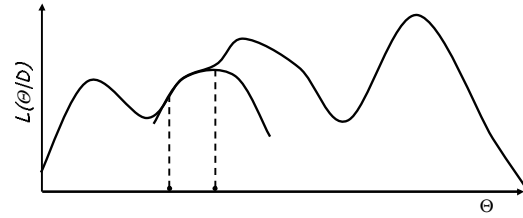


Example:

- ◆ We can rename the values of hidden variable H
- ◆ If H has two values, likelihood has two maxima
- ◆ In practice, many local maxima

73

EM: MLE from Incomplete Data



- ◆ Use current point to construct "nice" alternative function
- ◆ Max of new function scores \geq than current point

74

Expectation Maximization (EM)

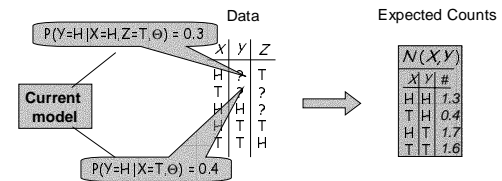
- ◆ A general purpose method for learning from incomplete data

Intuition:

- ◆ If we had true counts, we could estimate parameters
- ◆ But with missing values, counts are unknown
- ◆ We "complete" counts using probabilistic inference based on current parameter assignment
- ◆ We use completed counts as if real to re-estimate parameters

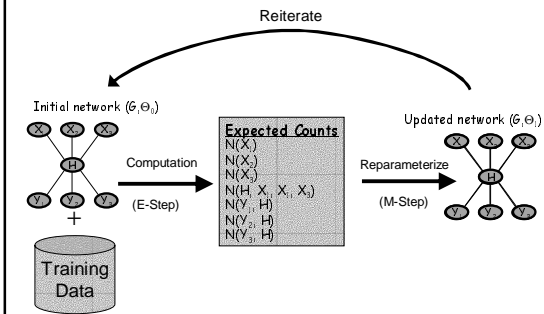
75

Expectation Maximization (EM)



76

Expectation Maximization (EM)



77

Expectation Maximization (EM)

Formal Guarantees:

- ◆ $L(\theta_1; D) \geq L(\theta_0; D)$
 - Each iteration improves the likelihood
- ◆ If $\theta_1 = \theta_0$, then θ_0 is a **stationary point** of $L(\theta; D)$
 - Usually, this means a local maximum

78

Expectation Maximization (EM)

Computational bottleneck:

- ◆ Computation of expected counts in E-Step
 - Need to compute posterior for each unobserved variable in each instance of training set
 - All posteriors for an instance can be derived from one pass of standard BN inference

79

Summary: Parameter Learning with Incomplete Data

- ◆ Incomplete data makes parameter estimation hard
- ◆ Likelihood function
 - Does not have closed form
 - Is multimodal
- ◆ Finding max likelihood parameters:
 - EM
 - Gradient ascent
- ◆ Both exploit inference procedures for Bayesian networks to compute expected sufficient statistics

80

Incomplete Data: Structure Scores

Recall, Bayesian score:

$$P(\mathcal{G} | D) \propto P(\mathcal{G})P(D | \mathcal{G})$$

$$= P(\mathcal{G}) \int P(D | \mathcal{G}, \theta)P(\theta | \mathcal{G})d\theta$$

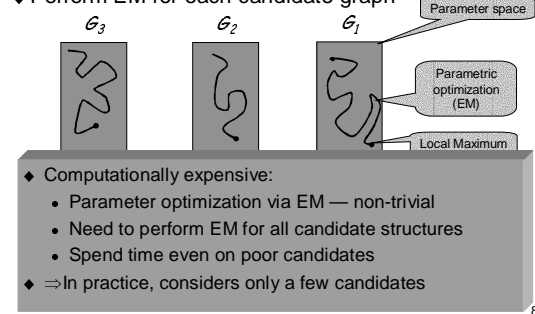
With incomplete data:

- ◆ Cannot evaluate marginal likelihood in closed form
- ◆ We have to resort to **approximations**:
 - Evaluate score around MAP parameters
 - Need to find MAP parameters (e.g., EM)

81

Naive Approach

- ◆ Perform EM for each candidate graph



- ◆ Computationally expensive:
 - Parameter optimization via EM — non-trivial
 - Need to perform EM for all candidate structures
 - Spend time even on poor candidates
- ◆ ⇒ In practice, considers only a few candidates

82

Structural EM

Recall, in complete data we had

- Decomposition ⇒ efficient search

Idea:

- ◆ Instead of optimizing the real score...
- ◆ Find **decomposable** alternative score
- ◆ Such that maximizing new score
 - ⇒ improvement in real score

83

Structural EM

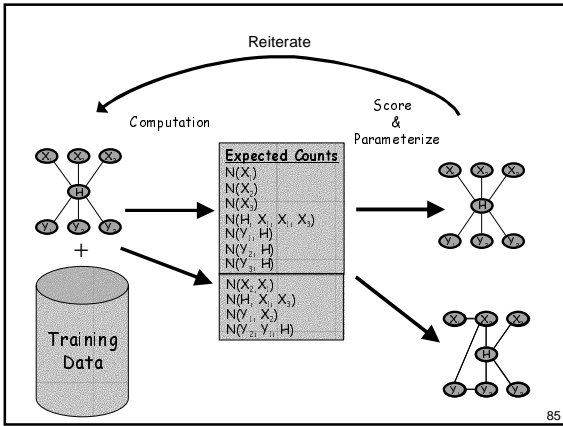
Idea:

- ◆ Use current model to help evaluate new structures

Outline:

- ◆ Perform search in (Structure, Parameters) space
- ◆ At each iteration, use current model for finding either:
 - Better scoring parameters: “parametric” EM step or
 - Better scoring structure: “structural” EM step

84



Example: Phylogenetic Reconstruction

Input: Biological sequences

Human CGTTGC...

Chimp CCTAGG...

Orang CGAACG...

....

Output: a phylogeny

10 billion years

leaf

An "instance" of evolutionary process

Assumption: positions are independent

86

Phylogenetic Model

- ◆ Topology: bifurcating
 - Observed species – $1 \dots N$
 - Ancestral species – $N+1 \dots 2N-2$
- ◆ Lengths $t = \{t_{i,j}\}$ for each branch (i,j)
- ◆ Evolutionary model:
 - $P(A \text{ changes to } T | 10 \text{ billion yrs})$

87

Phylogenetic Tree as a Bayes Net

- ◆ Variables: Letter at each position for each species
 - Current day species – observed
 - Ancestral species - hidden
- ◆ BN Structure: Tree topology
- ◆ BN Parameters: Branch lengths (time spans)

Main problem: Learn topology

If ancestral were observed
 \Rightarrow easy learning problem (learning trees)

88

Algorithm Outline

- \rightarrow Compute expected pairwise stats
- \rightarrow Weights: Branch scores

Original Tree (T^0, t^0)

89

Algorithm Outline

- \rightarrow Compute expected pairwise stats
- \rightarrow Weights: Branch scores
- \rightarrow Find: $T^1 = \text{argmax}_T \sum_{(i,j) \in T} w_{i,j}$

Pairwise weights

$O(N^2)$ pairwise statistics suffice to evaluate all trees

90

Algorithm Outline

- Compute expected pairwise stats
- Weights: Branch scores
- Find: $T' = \operatorname{argmax}_T \sum_{(i,j) \in T} w_{i,j}$
- Construct bifurcation T_I

Max. Spanning Tree

91

Algorithm Outline

- Compute expected pairwise stats
- Weights: Branch scores
- Find: $T' = \operatorname{argmax}_T \sum_{(i,j) \in T} w_{i,j}$
- Construct bifurcation T_I
- Theorem: $L(T_I, t_I) \geq L(T_0, t_0)$

New Tree

Repeat until convergence...

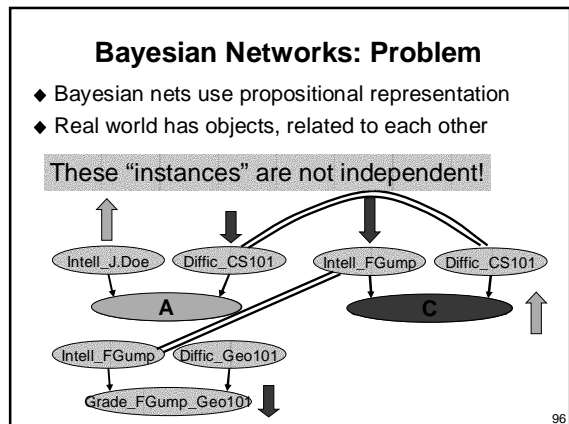
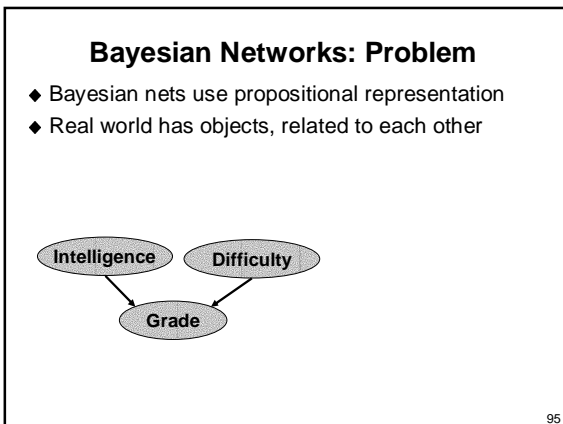
92

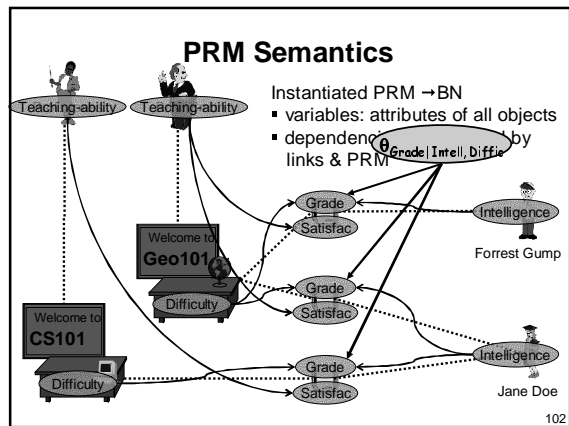
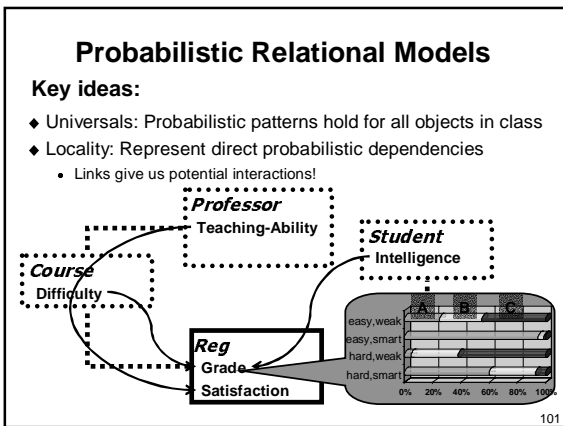
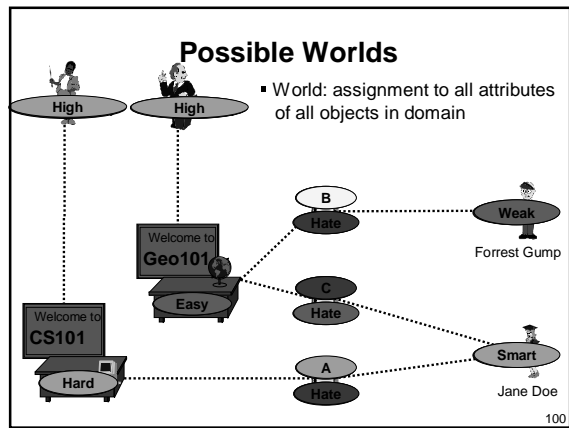
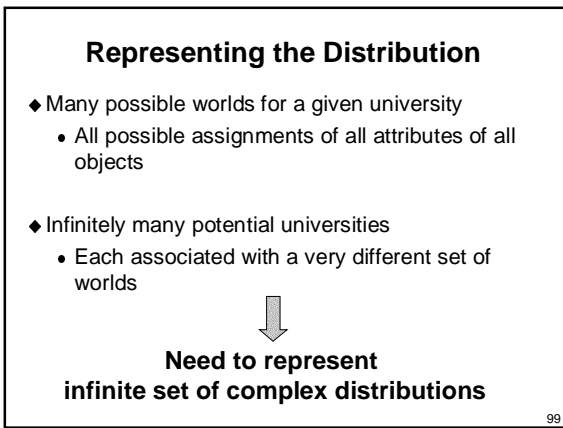
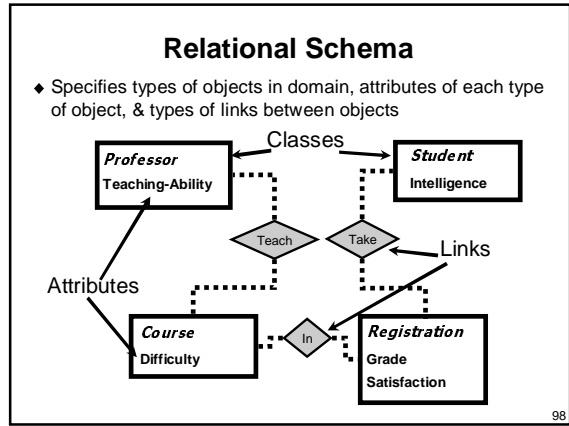
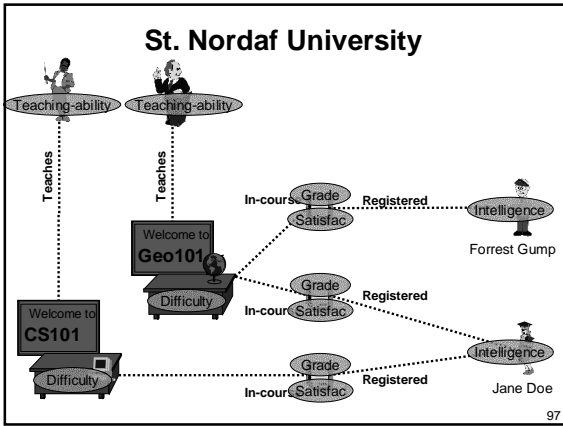
Real Life Data

	Lysozyme c	Mitochondrial genomes
# sequences	43	34
# pos	122	3,578
Log-likelihood	Traditional approach	-2,916.2
	Structural EM Approach	-2,892.1
	Difference per position	0.19
		1.03

93

- ### Overview
- ◆ Introduction
 - ◆ Parameter Estimation
 - ◆ Model Selection
 - ◆ Structure Discovery
 - ◆ Incomplete Data
 - ◆ **Learning from Structured Data**
- 94





The Web of Influence

- Objects are all correlated
- Need to perform inference over entire model
- For large databases, use approximate inference:
 - Loopy belief propagation

103

PRM Learning: Complete Data

- Introduce prior over parameters
- Update prior with sufficient statistics:

$$\text{Count}(\text{Reg. Grade}=A, \text{Reg. Course. Diff}=lo, \text{Reg. Student. Intel}=hi)$$

104

PRM Learning: Incomplete Data

- Use expected sufficient statistics
- But, everything is correlated:
 - E-step uses (approx) inference over entire model

105

A Web of Data

[Craven et al.]

106

Standard Approach

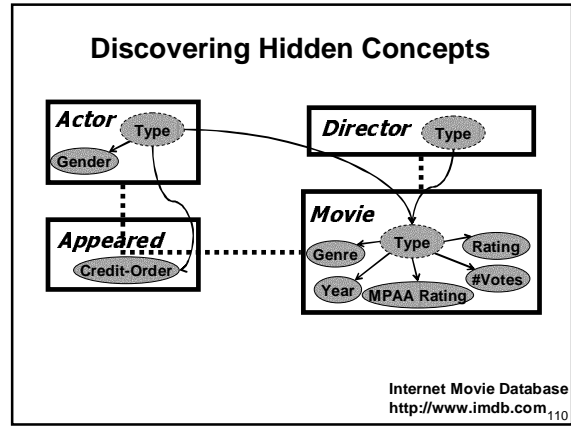
107

What's in a Link

108

Discovering Hidden Concepts

Internet Movie Database
<http://www.imdb.com>



Web of Influence, Yet Again

Movies	Actors	Directors
Wizard of Oz Cinderella Sound of Music The Love Bug Pollyanna The Parent Trap Mary Poppins Swiss Family Robinson Terminator 2 Batman Batman Forever Mission: Impossible GoldenEye Starship Troopers Hunt for Red October ...	Sylvester Stallone Bruce Willis Harrison Ford Steven Seagal Kurt Russell Kevin Costner Jean-Claude Van Damme Arnold Schwarzenegger Anthony Hopkins Robert De Niro Tommy Lee Jones Harvey Keitel Morgan Freeman Gary Oldman ...	Alfred Hitchcock Stanley Kubrick David Lean Milos Forman Terry Gilliam Francis Coppola Steven Spielberg Tim Burton Tony Scott James Cameron John McTiernan Joel Schumacher ...

- ### Conclusion
- ◆ Many distributions have combinatorial dependency structure
 - ◆ Utilizing this structure is good
 - ◆ Discovering this structure has implications:
 - To density estimation
 - To knowledge discovery
 - ◆ Many applications
 - Medicine
 - Biology
 - Web

The END

Thanks to

- ◆ Gal Elidan
- ◆ Dana Pe'er
- ◆ Lise Getoor
- ◆ Eran Segal
- ◆ Moises Goldszmidt
- ◆ Ben Taskar
- ◆ Matan Ninio

Slides will be available from:
<http://www.cs.huji.ac.il/~nir/>
<http://robotics.stanford.edu/~koller/>