

**CS 2750 Machine Learning  
Lecture 17**

**Density estimation with hidden  
variables and missing values**

Milos Hauskrecht  
[milos@cs.pitt.edu](mailto:milos@cs.pitt.edu)  
5329 Sennott Square

---

CS 2750 Machine Learning

**Administration**

**Midterm: Wednesday, March 17, 2004**

- **In class**
- **Closed book**
- **Material covered by Spring break, excluding learning of the BBN structures**
- **Last year midterm is on the web**

**No new homework**

---

CS 2750 Machine Learning

## Project proposals

**Due: Wednesday, March 24, 2004**

- **1-2 pages long**

### **Proposal**

- **Written proposal:**
  1. Outline of a learning problem, type of data you have available. Why is the problem important?
  2. Learning methods you plan to try and implement for the problem. References to previous work.
  3. How do you plan to test, compare learning approaches
  4. Schedule of work (approximate timeline of work)
- **A PPT (3 slide) summary of points 1-4**

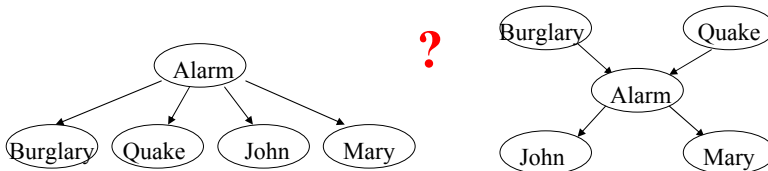
## Learning the structure of the BBN

## Model selection

- **BBN has two components:**
  - **Structure of the network** (models conditional independences)
  - **A set of parameters** (conditional child-parent distributions)

We already know how to learn the parameters for the fixed structure

But how to learn the structure of the BBN?



CS 2750 Machine Learning

## Learning the structure

Criteria we can choose to score the structure  $S$

- **Marginal likelihood**

$$\text{maximize } P(D | S, \xi)$$

$\xi$  - represents the prior knowledge

- **Maximum posterior probability**

$$\text{maximize } P(S | D, \xi)$$

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

How to compute marginal likelihood  $P(D | S, \xi)$  ?

CS 2750 Machine Learning

## Learning of BBNs

- **Notation:**

- $i$  ranges over all possible variables  $i=1, \dots, n$
- $j=1, \dots, q$  ranges over all possible parent combinations
- $k=1, \dots, r$  ranges over all possible variable values
- $\Theta$  - parameters of the BBN

$\Theta_{ij}$  is a vector of  $\Theta_{ijk}$  representing parameters of the conditional probability distribution; such that  $\sum_{k=1}^r \Theta_{ijk} = 1$

$N_{ijk}$  - a number of instances in the dataset where parents of variable  $X_i$  take on values  $j$  and  $X_i$  has value  $k$

$$N_{ij} = \sum_{k=1}^r N_{ijk}$$

$\alpha_{ijk}$  - prior counts (parameters of Beta and Dirichlet priors)

$$\alpha_{ij} = \sum_{k=1}^r \alpha_{ijk}$$

## Marginal likelihood

- Integrate over all possible parameter settings

$$P(D | S, \xi) = \int_{\Theta} P(D | S, \Theta, \xi) p(\Theta | S, \xi) d\Theta$$

- Using the assumption of parameter and sample independence

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

- We can use **log-likelihood score** instead

$$\log P(D | S, \xi) = \sum_{i=1}^n \left\{ \sum_{j=1}^{q_i} \log \frac{\Gamma(\alpha_{ij})}{\Gamma(\alpha_{ij} + N_{ij})} + \sum_{k=1}^{r_i} \log \frac{\Gamma(\alpha_{ijk} + N_{ijk})}{\Gamma(\alpha_{ijk})} \right\}$$

**Score is decomposable along variables !!!**

## Computing the marginal likelihood

- From the iid assumption:

$$P(\mathbf{D} | S, \Theta) = \prod_{h=1}^N \prod_{i=1}^n P(x_i^h | \text{parents}_i^h, \Theta)$$

- Let  $r_i$  = number of values that attribute  $x_i$  can take  
 $q_i$  = number of possible parent combinations  
 $N_{ijk}$  = number of cases in  $\mathbf{D}$  where  $x_i$  has value  $k$  and parents with values  $j$ .

$$\begin{aligned} &= \prod_i^n \prod_j^{q_i} \prod_k^{r_i} P(x_i = k | \text{parents}_i = j, \Theta)^{N_{ijk}} \\ &= \prod_i^n \prod_j^{q_i} \prod_k^{r_i} \theta_{ijk}^{N_{ijk}} \end{aligned}$$

## Computing the marginal likelihood

- From parameter independence

$$p(\Theta | S, \xi) = \prod_{i=1}^n \prod_{j=1}^{q_i} p(\Theta_{ij} | S, \xi)$$

- Priors for  $p(\Theta_{ij} | S, \xi)$

- $\Theta_{ij} = (\Theta_{ij1}, \dots, \Theta_{ijr_i})$  is a vector of parameters;
- we use a Dirichlet distribution with parameters  $\alpha$  to represent it

$$P(\Theta_{ij} | S, \xi) = P(\Theta_{ij1}, \dots, \Theta_{ijr_i} | S, \xi) = \text{Dirichlet}(\Theta_{ij1}, \dots, \Theta_{ijr_i} | \alpha)$$

$$\begin{aligned} &= \frac{\Gamma(\sum_{k=1}^{r_i} \alpha_{ijk})}{\prod_{k=1}^{r_i} \Gamma(\alpha_{ijk})} \prod_{k=1}^{r_i} \Theta_{ijk}^{\alpha_{ijk}-1} \end{aligned}$$

## Computing the marginal likelihood

- **Combine things together:**

$$\begin{aligned}
 P(D | S_i) &= \int_{\Theta} P(D | S_i, \Theta) P(\Theta | S_i) d\Theta \\
 &= \int \prod_i^n \prod_j^{q_i} \prod_k^{r_i} \Theta_{ijk}^{N_{ijk}} \cdot \frac{\Gamma(\sum_{k=1}^{r_i} \alpha_{ijk})}{\prod_{k=1}^{r_i} \Gamma(\alpha_{ijk})} \prod_{k=1}^{r_i} \Theta_{ijk}^{\alpha_{ijk}-1} d\Theta \\
 &= \prod_i^n \prod_j^{q_i} \frac{\Gamma(\sum_{k=1}^{r_i} \alpha_{ijk})}{\prod_{k=1}^{r_i} \Gamma(\alpha_{ijk})} \int \prod_{k=1}^{r_i} \Theta_{ijk}^{N_{ijk} + \alpha_{ijk} - 1} d\Theta \\
 &= \prod_i^n \prod_j^{q_i} \frac{\Gamma(\alpha_{ij})}{\prod_{k=1}^{r_i} \Gamma(\alpha_{ijk})} \cdot \frac{\prod_{k=1}^{r_i} \Gamma(\alpha_{ijk} + N_{ijk})}{\Gamma(\alpha_{ij} + N_{ij})}
 \end{aligned}$$

CS 2750 Machine Learning

## Learning the structure

- **Likelihood of data for the BBN** (structure + parameters)

$$P(D | S, \Theta, \xi)$$

measures the goodness of fit of the BBN to data

- **Marginal likelihood** (for the structure only)

$$P(D | S, \xi)$$

- **Does not measure only a goodness of fit. It is:**
  - different for structures of different complexity
  - Incorporates preferences towards simpler structures, implements **Occam's razor** !!!!

CS 2750 Machine Learning

## Occam's Razor

- Why there is a preference towards simpler structures ?

Rewrite marginal likelihood as

$$P(D | S, \xi) = \frac{\int_{\Theta} P(D | S, \Theta, \xi) p(\Theta | S, \xi) d\Theta}{\int_{\Theta} p(\Theta | S, \xi) d\Theta}$$

We know that  $\int_{\Theta} p(\Theta | S, \xi) d\Theta = 1$

**Interpretation:** in more complex structures there are more ways how parameters can be set badly

- **The numerator:** count of good assignments
- **The denominator:** count of all assignments

## Approximations of probabilistic scores

Approximations of the marginal likelihood and posterior scores

- **Information based measures**

- Akaike criterion
- Bayesian information criterion (BIC)
- Minimum description length (MDL)

- Reflect the tradeoff between the fit to data and preference towards simpler structures

Example: **Akaike criterion.**

**Maximize:**  $score(S) = \log P(D | S, \Theta_{ML}, \xi) - \text{compl}(S)$

**Bayesian information criterion (BIC)**

**Maximize:**  $score(S) = \log P(D | S, \Theta_{ML}, \xi) - \frac{1}{2} \text{compl}(S) \log N$

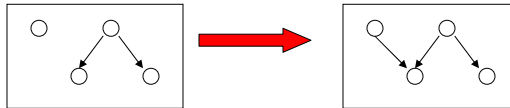
## Optimizing the structure

Finding the best structure is a **combinatorial optimization** problem

- A good feature: the score **is decomposable along variables**:

$$\log P(D | S, \xi) = \sum_{i=1}^n \left\{ \sum_{j=1}^{q_i} \log \frac{\Gamma(\alpha_{ij})}{\Gamma(\alpha_{ij} + N_{ij})} + \sum_{k=1}^{r_i} \log \frac{\Gamma(\alpha_{ijk} + N_{ijk})}{\Gamma(\alpha_{ijk})} \right\}$$

**Algorithm idea:** Search the space of structures using local changes (additions and deletions of a link)



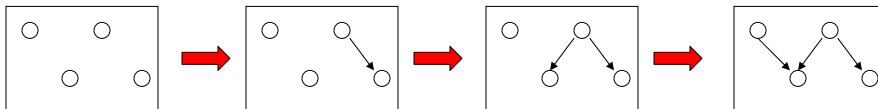
**Advantage:**

- we do not have to compute the whole score from scratch
- Recompute the partial score for the affected variable

## Optimizing the structure. Algorithms

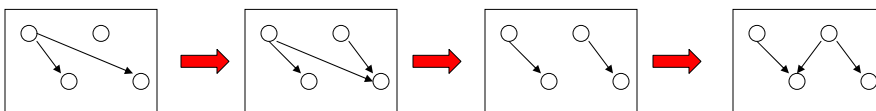
- **Greedy search**

- Start from structure with no links
- Add a link that yields the best score improvement



- **Metropolis algorithm (with simulated annealing)**

- Local additions and deletions
- Avoids being trapped in “local” optimal





## Density estimation with hidden variables and missing values

CS 2750 Machine Learning

## Learning probability distribution

### Basic learning settings:

- A set of random variables  $\mathbf{X} = \{X_1, X_2, \dots, X_n\}$
- **A model of the distribution** over variables in  $\mathbf{X}$  with parameters  $\Theta$
- **Data**  $D = \{D_1, D_2, \dots, D_N\}$   
**s.t.**  $D_i = (x_1^i, x_2^i, \dots, x_n^i)$

**Objective:** find parameters  $\hat{\Theta}$  that describe the data

### Assumptions considered so far:

- Known parameterizations
- No hidden variables
- No-missing values

CS 2750 Machine Learning

## Hidden variables

### Modeling assumption:

Variables  $\mathbf{X} = \{X_1, X_2, \dots, X_n\}$  are related through hidden variables

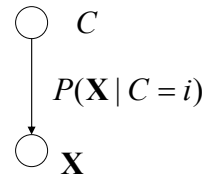
### Why to add hidden variables?

- More flexibility in describing the distribution  $P(\mathbf{X})$
- Smaller parameterization of  $P(\mathbf{X})$ 
  - New independences can be introduced via hidden variables

### Example:

- Latent variable models
  - hidden classes (categories)

Hidden class variable

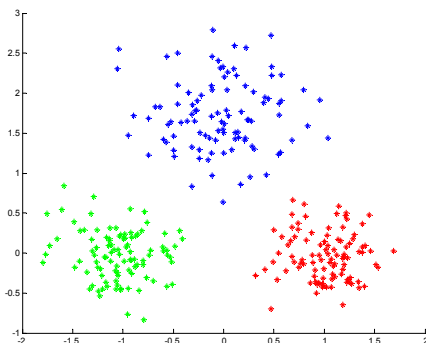


CS 2750 Machine Learning

## Hidden variable model. Example.

- We want to represent the probability model of a population in a two dimensional space  $\mathbf{X} = \{X_1, X_2\}$

### Observed data

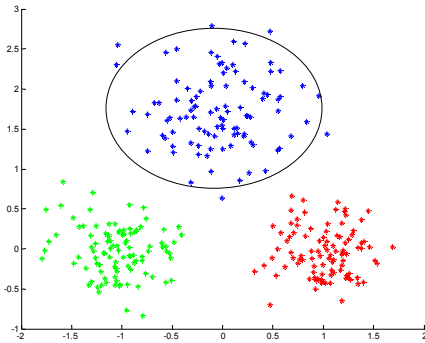


CS 2750 Machine Learning

## Hidden variable model

- We want to represent the probability model of a population in a two dimensional space  $\mathbf{X} = \{X_1, X_2\}$

### Observed data

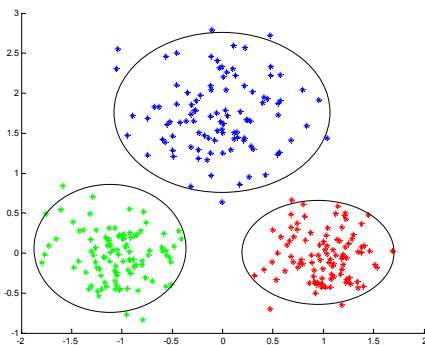


CS 2750 Machine Learning

## Hidden variable model

- We want to represent the probability model of a population in a two dimensional space  $\mathbf{X} = \{X_1, X_2\}$

### Observed data

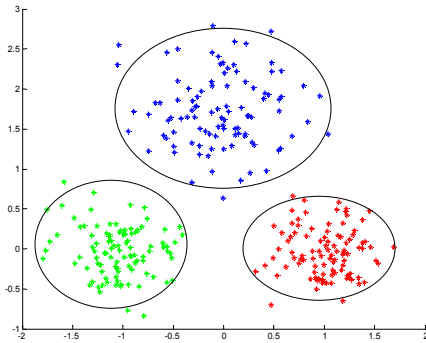


CS 2750 Machine Learning

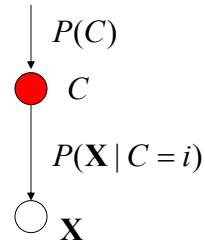
## Hidden variable model

- We want to represent the probability model of a population in a two dimensional space  $\mathbf{X} = \{X_1, X_2\}$

**Observed data**



**Model** : 3 Gaussians with a hidden class variable



CS 2750 Machine Learning

## Mixture of Gaussians

Probability of the occurrence of a data point  $\mathbf{x}$  is modeled as

$$p(\mathbf{x}) = \sum_{i=1}^k p(C = i) p(\mathbf{x} | C = i)$$

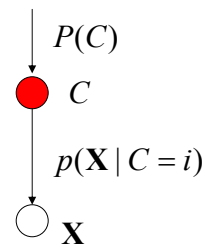
where

$$p(C = i)$$

= probability of a data point coming from class  $C=i$

$$p(\mathbf{x} | C = i) \approx N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$

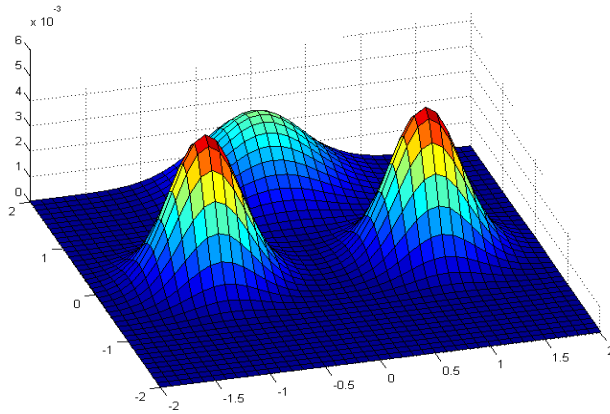
= class-conditional density (modeled as Gaussian) for class  $i$



CS 2750 Machine Learning

## Mixture of Gaussians

- Density function for the Mixture of Gaussians model



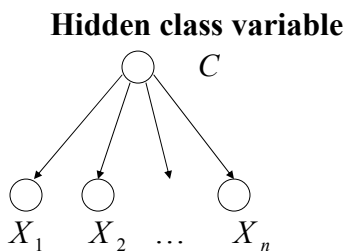
CS 2750 Machine Learning

## Naïve Bayes with a hidden class variable

**Introduction of a hidden variable can reduce the number of parameters defining  $P(\mathbf{X})$**

**Example:**

- Naïve Bayes model with a hidden class variable



Attributes are independent given the class

- **Useful in customer profiles**
  - Class value = type of customers

CS 2750 Machine Learning

## Missing values

A set of random variables  $\mathbf{X} = \{X_1, X_2, \dots, X_n\}$

• **Data**  $D = \{D_1, D_2, \dots, D_N\}$

• **But some values are missing**

$$D_i = (x_1^i, x_3^i, \dots, x_n^i)$$

Missing value of  $x_2^i$

$$D_{i+1} = (x_3^i, \dots, x_n^i)$$

Missing values of  $x_1^i, x_2^i$

Etc.

• **Example: medical records**

• **We still want to estimate parameters of**  $P(\mathbf{X})$

## Density estimation

**Goal: Find the set of parameters**  $\hat{\Theta}$

**Estimation criteria:**

- **ML**  $\max_{\Theta} p(D | \Theta, \xi)$
- **Bayesian**  $p(\Theta | D, \xi)$

**Possible optimization methods for the ML:** gradient-ascent, conjugate gradient, Newton-Rhapon, etc.

• **Problem:** No or very small advantage from the structure of the corresponding belief network

**Expectation-maximization (EM) method**

- An alternative optimization method
- Suitable when there are missing or hidden values
- **Takes advantage of the structure of the belief network**