

# CSE 473: Artificial Intelligence

## Bayesian Networks - Learning

Dieter Fox

Slides adapted from Dan Weld, Jack Breese, Dan Klein, Daphne Koller, Stuart Russell, Andrew Moore & Luke Zettlemoyer

# Space of ML Problems

Type of Supervision  
(eg, Experience, Feedback)

	Labeled Examples	Reward	Nothing
Discrete Function	Classification		Clustering
Continuous Function	Regression		
Policy	Apprenticeship Learning	Reinforcement Learning	

What is Being Learned?

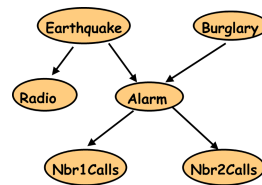
2

# Learning Topics

- Learning Parameters for a Bayesian Network
  - Fully observable
    - Maximum Likelihood (ML)
    - Maximum A Posteriori (MAP)
    - Bayesian
  - Hidden variables (EM algorithm)
- Learning Structure of Bayesian Networks

© Daniel S. Weld

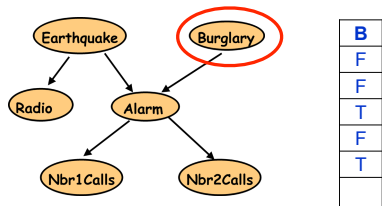
# Parameter Estimation and Bayesian Networks



E	B	R	A	J	M
T	F	T	T	F	T
F	F	F	F	F	T
F	T	F	T	T	T
F	F	F	T	T	T
F	T	F	F	F	F
...					

- We have:
- Bayes Net **structure** and **observations**
  - We need: Bayes Net **parameters**

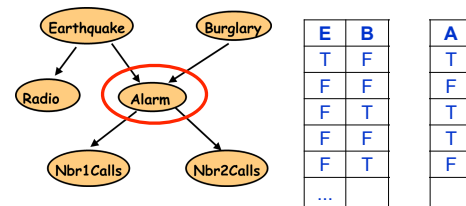
# Parameter Estimation and Bayesian Networks



$$P(B) = ? = 0.4$$

$$P(\neg B) = 1 - P(B) = 0.6$$

# Parameter Estimation and Bayesian Networks



$$P(A|E,B) = ?$$

$$P(A|E,\neg B) = ?$$

$$P(A|\neg E,B) = ?$$

$$P(A|\neg E,\neg B) = ?$$

### Parameter Estimation and Bayesian Networks

E	B
T	F
F	F
F	T
F	F
F	T
...	

A
T
F
T
T
F

$P(A|E,B) = ?$   
 $P(A|E,\neg B) = ?$   
 $P(A|\neg E,B) = ?$   
 $P(A|\neg E,\neg B) = 0.5$

### Parameter Estimation and Bayesian Networks

B
F
F
T
F
T

$P(B) =$   $+ \text{data} =$

Now compute either MAP or Bayesian estimate

### Parameter Estimation and Bayesian Networks

B
F
F
T
F
T

**Prior**  
 $P(B|\text{data}) = \text{Beta}(1,4)$  “+ data” =  $(3,7)$

B	$\neg B$
.3	.7

Prior  $P(B) = 1/(1+4) = 20\%$  with equivalent sample size 5

### Parameter Estimation and Bayesian Networks

E	B
T	F
F	F
F	T
F	F
F	T
...	

A
T
F
T
T
F

$P(A|E,B) = ?$   
 $P(A|E,\neg B) = ?$   
 $P(A|\neg E,B) = ?$   
 $P(A|\neg E,\neg B) = ?$

### Parameter Estimation and Bayesian Networks

E	B
T	F
F	F
F	T
F	F
F	T
...	

A
T
F
T
T
F

**Prior**  
 $P(A|E,B) = ?$   
 $P(A|E,\neg B) = ?$   
 $P(A|\neg E,B) = ?$  **Beta(2,3)**  
 $P(A|\neg E,\neg B) = ?$

### Parameter Estimation and Bayesian Networks

E	B
T	F
F	F
F	T
F	F
F	T
...	

A
T
F
T
T
F

**Prior**  
 $P(A|E,B) = ?$   
 $P(A|E,\neg B) = ?$   
 $P(A|\neg E,B) = ?$  **Beta(2,3)** + data =  $(3,4)$   
 $P(A|\neg E,\neg B) = ?$

What if we *don't* know structure?

### Learning The Structure of Bayesian Networks

- Search through the space...
  - of possible network structures!
  - (for now, assume we observe all variables)
- For each structure, learn parameters
- Pick the one that fits observed data best
  - Caveat – won't we end up fully connected????

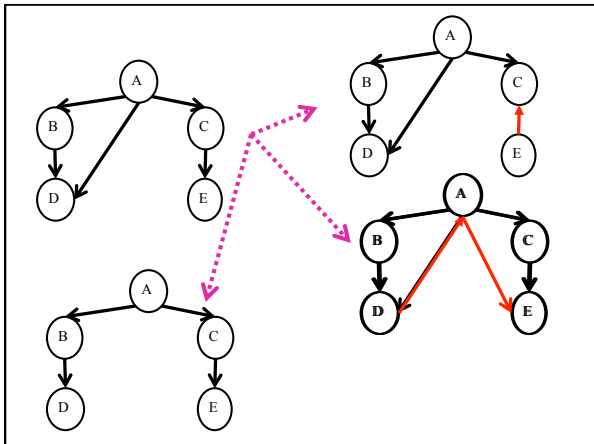
When scoring, add a penalty  
model complexity

### Learning The Structure of Bayesian Networks

- Search through the space
- For each structure, learn parameters
- Pick the one that fits observed data best
  - Penalize complex models
- Problem?
  - Exponential number of networks!
  - And we need to learn parameters for each!
  - Exhaustive search out of the question!

### Structure Learning as Search

- Local Search
  - Start with some network structure
  - Try to make a change (add or delete or reverse edge)
  - See if the new network is any better
- What should the initial state be?
  - Uniform prior over random networks?
  - Based on prior knowledge?
  - Empty network?
- How do we evaluate networks?



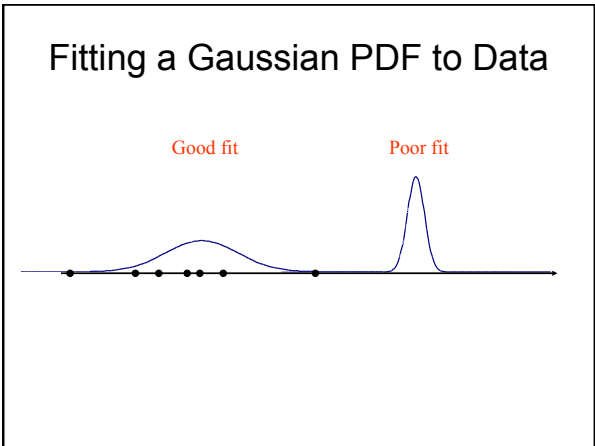
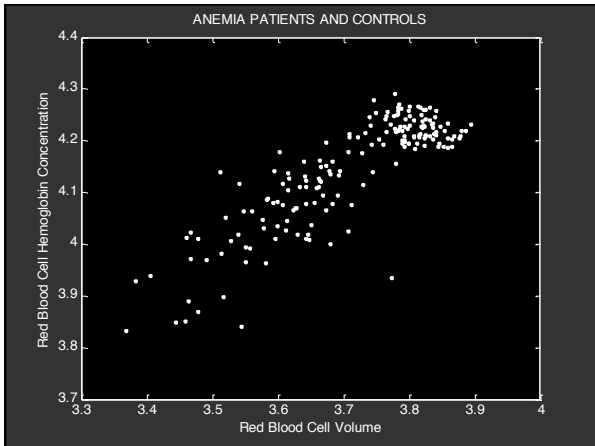
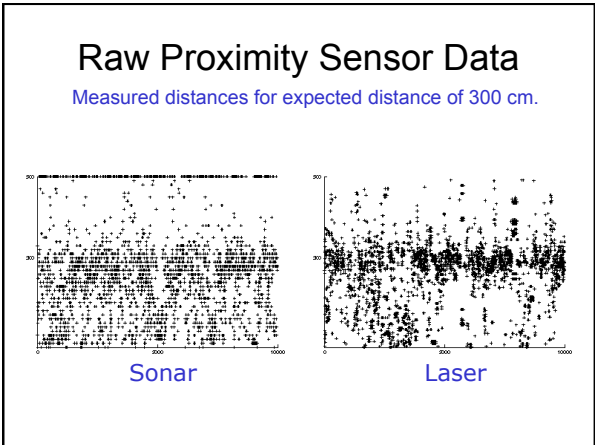
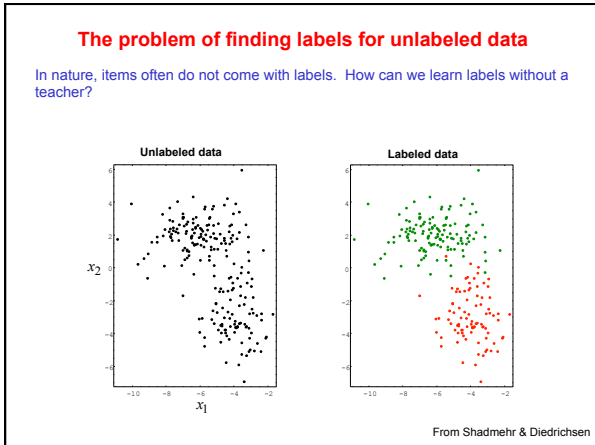
### Scoring a Bayes Net Structure

- Bayesian Information Criterion (BIC)
  - $P(D | BN)$  – penalty
  - Penalty =  $\frac{1}{2} (\# \text{ parameters}) \text{ Log} (\# \text{ data points})$

# Expectation Maximization and Gaussian Mixtures

CSE 473

- ## Feedback in Learning
- Supervised learning: correct answers for each example
  - Unsupervised learning: correct answers not given
  - Reinforcement learning: occasional rewards



### Fitting a Gaussian PDF to Data

- Suppose  $y = y_1, \dots, y_n, \dots, y_N$  is a set of  $N$  data values
- Given a Gaussian PDF  $p$  with mean  $\mu$  and std dev  $\sigma$ , define:

$$p(y | \mu, \sigma) = \prod_{n=1}^N p(y_n | \mu, \sigma) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \frac{(y_n - \mu)^2}{\sigma^2}}$$

- How do we choose  $\mu$  and  $\sigma$  to maximise this probability?

Fisher, 1922

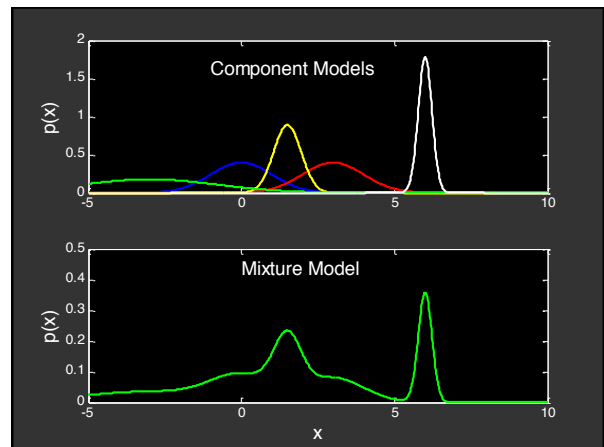
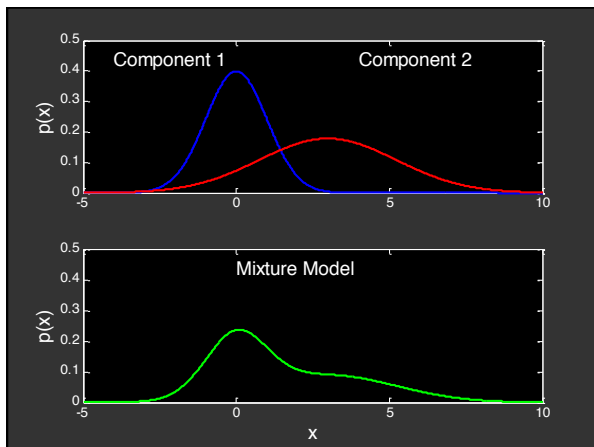
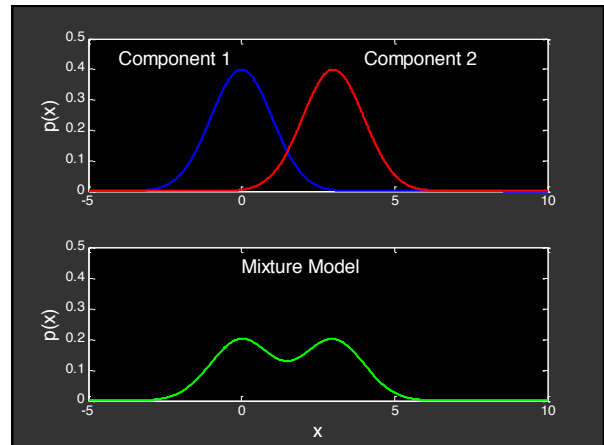
### Maximum Likelihood Estimation

- Define the best fitting Gaussian to be the one such that  $p(y|\mu, \sigma)$  is maximized.
- Terminology:
  - $p(y | \mu, \sigma)$ , thought of as a function of  $y$  is the **probability (density)** of  $y$
  - $p(y | \mu, \sigma)$ , thought of as a function of  $\mu, \sigma$  is the **likelihood** of  $\mu, \sigma$
- Maximizing  $p(y | \mu, \sigma)$  with respect to  $\mu, \sigma$  is called **Maximum Likelihood (ML)** estimation of  $\mu, \sigma$

### ML estimation of $\mu, \sigma$

- Intuitively:
  - The maximum likelihood estimate of  $\mu$  should be the average value of  $y_1, \dots, y_N$ , (the sample mean)
  - The maximum likelihood estimate of  $\sigma$  should be the variance of  $y_1, \dots, y_N$ . (the sample variance)
- This turns out to be true:  $p(y | \mu, \sigma)$  is maximized by setting:

$$\mu = \frac{1}{N} \sum_{n=1}^N y_n, \quad \sigma = \frac{1}{N} \sum_{n=1}^N (y_n - \mu)^2$$



### Mixtures

If our data is not labeled, we can hypothesize that:

- There are exactly  $m$  classes in the data:  $y \in \{1, 2, \dots, m\}$
- Each class  $y$  occurs with a specific frequency:  $P(y)$
- Examples of class  $y$  are governed by a specific distribution:  $p(x|y)$

According to our hypothesis, each example  $x^{(i)}$  must have been generated from a specific "mixture" distribution:

$$p(x) = \sum_{j=1}^m P(y=j) p(x|y=j)$$

We might hypothesize that the distributions are Gaussian:

Parameters of the distributions  $\theta = \{P(y=1), \mu_1, \Sigma_1, \dots, P(y=m), \mu_m, \Sigma_m\}$

$$p(x|\theta) = \sum_{j=1}^m P(y=j) \mathcal{N}(x|\mu_j, \Sigma_j)$$

Mixing proportions
Normal distribution

### Graphical Representation of Gaussian Mixtures

$$p(x) = \sum_{i=1}^3 P(y=i) p(x|y=i, \mu_i, \sigma_i)$$

## Learning of mixture models

### Learning Mixtures from Data

Consider fixed  $K = 2$

e.g., unknown parameters  $Q = \{m_1, s_1, m_2, s_2, a_1\}$

Given data  $D = \{x_1, \dots, x_N\}$ , we want to find the parameters  $Q$  that "best fit" the data

### 1977: The EM Algorithm

- **Dempster, Laird, and Rubin**
  - General framework for likelihood-based parameter estimation with missing data
    - start with initial guesses of parameters
    - E-step: estimate memberships given params
    - M-step: estimate params given memberships
    - Repeat until convergence
  - Converges to a **local** maximum of likelihood
  - E-step and M-step are often computationally simple
  - Can incorporate priors over parameters

### EM for Mixture of Gaussians

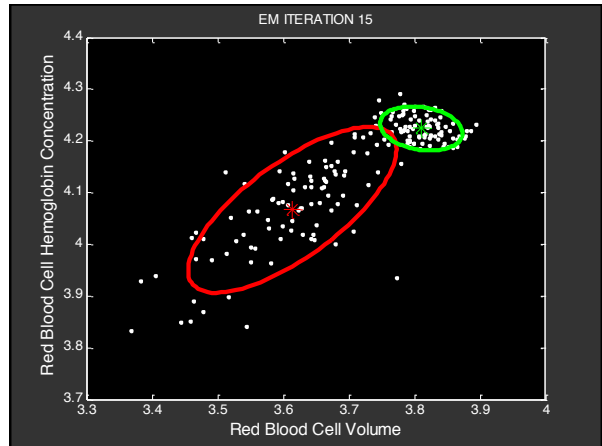
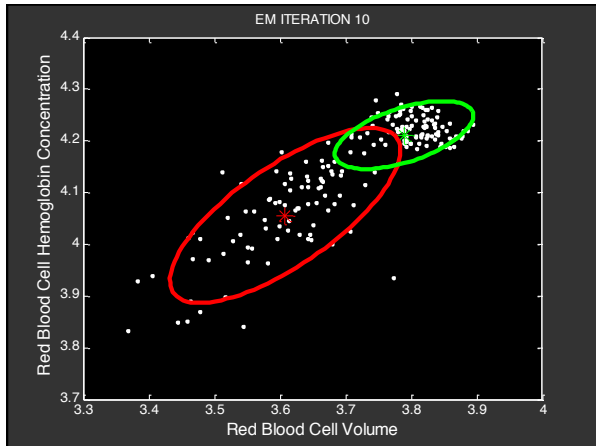
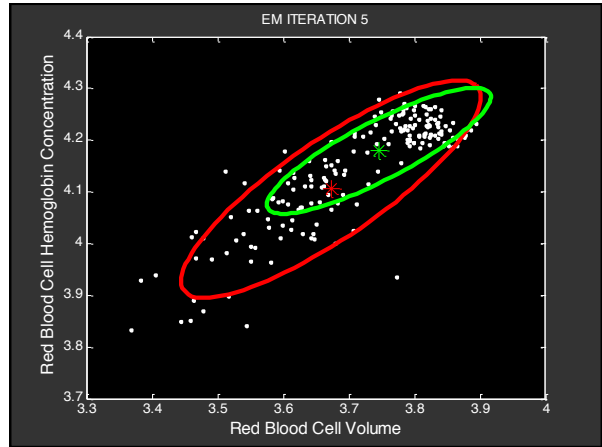
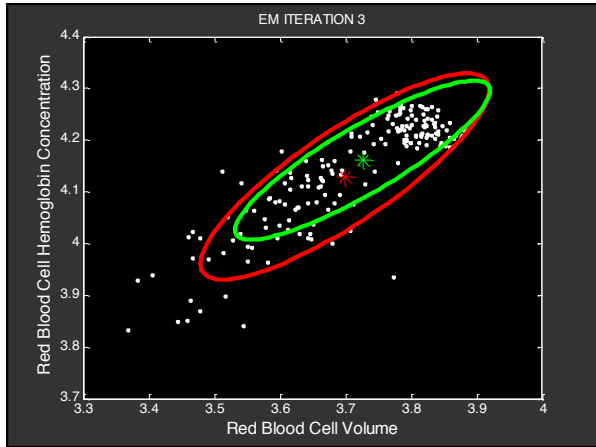
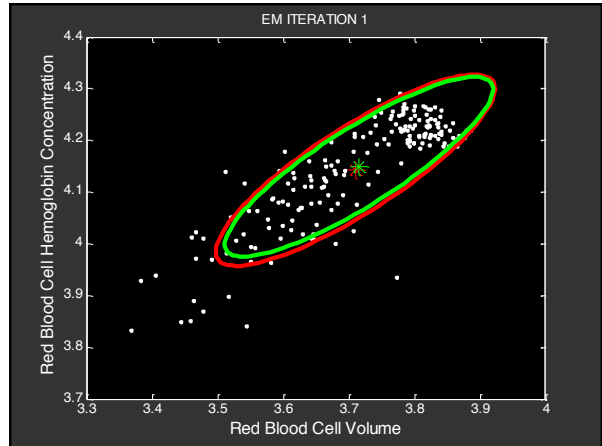
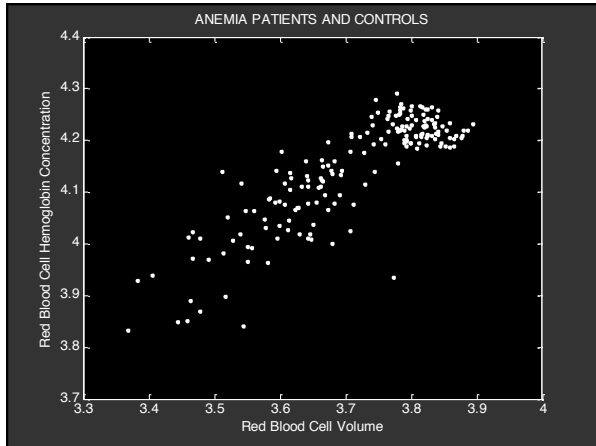
- **E-step:** Compute probability that point  $x_j$  was generated by component  $i$ :
 
$$p_{ij} = \alpha P(x_j | C=i) P(C=i)$$

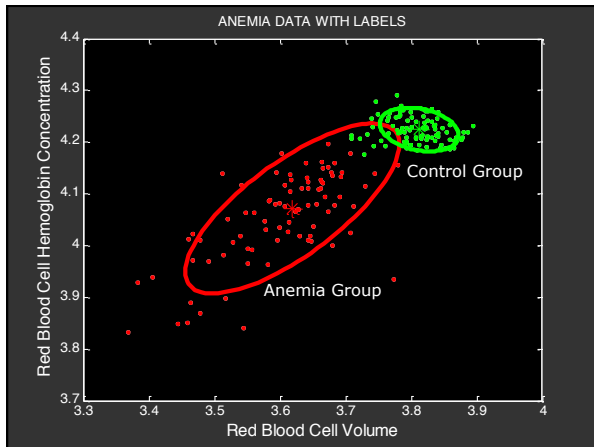
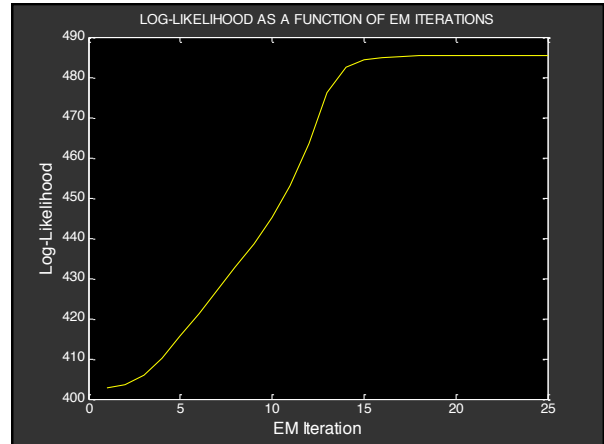
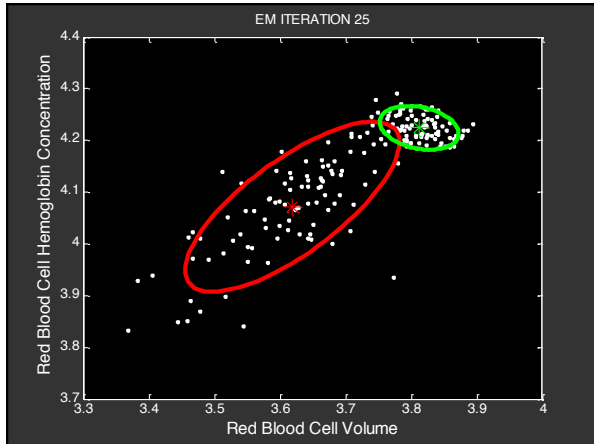
$$P_i = \sum_j p_{ij}$$
- **M-step:** Compute new mean, covariance, and component weights:
 
$$\mu_i \leftarrow \sum_j p_{ij} x_j / P_i$$

$$\sigma^2 \leftarrow \sum_j p_{ij} (x_j - \mu_i)^2 / P_i$$

$$w_i \leftarrow P_i$$

© D. Weld and D. Fox  
36





### Mixture Density

$$P(z | x, m) = \begin{pmatrix} \alpha_{\text{init}} \\ \alpha_{\text{unexp}} \\ \alpha_{\text{max}} \\ \alpha_{\text{rand}} \end{pmatrix}^T \begin{pmatrix} P_{\text{init}}(z | x, m) \\ P_{\text{unexp}}(z | x, m) \\ P_{\text{max}}(z | x, m) \\ P_{\text{rand}}(z | x, m) \end{pmatrix}$$

How can we determine the model parameters?

