

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

What is Being Learned?	Type of Supervision (eg, Experience, Feedback)		
	Labeled Examples	Reward	Nothing
	Discrete Function	Classification	Clustering
	Continuous Function	Regression	
Policy	Apprenticeship Learning	Reinforcement Learning	

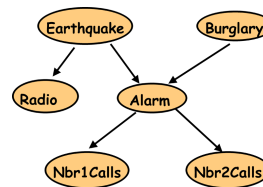
2

Learning Topics

- Learning Parameters for a Bayesian Network
 - Fully observable
 - Hidden variables (EM algorithm)
- Learning Structure of Bayesian Networks

© Daniel S. Weld

Parameter Estimation and Bayesian Networks

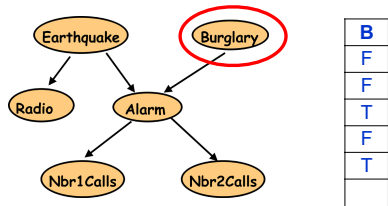


E	B	R	A	J	M
T	F	T	T	F	T
F	F	F	F	F	T
F	T	F	T	T	T
F	F	F	T	T	T
F	T	F	F	F	F
...					

We have:

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

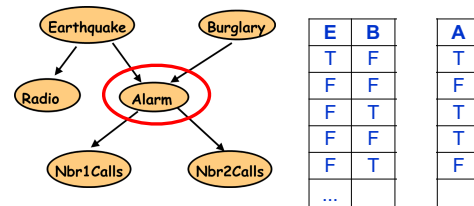
Parameter Estimation and Bayesian Networks



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

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

Parameter Estimation and Bayesian Networks



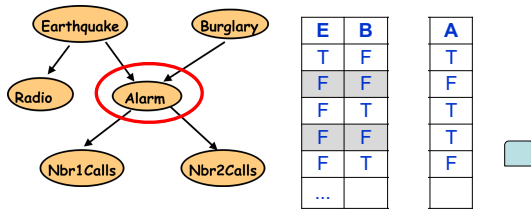
$$P(A|E,B) = ?$$

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

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

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

Parameter Estimation and Bayesian Networks



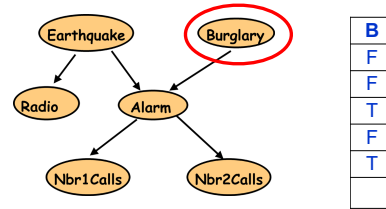
$$P(A|E,B) = ?$$

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

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

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

Parameter Estimation and Bayesian Networks

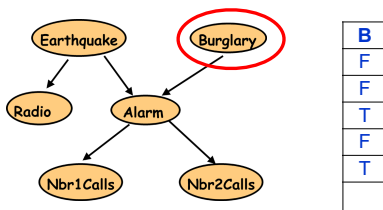


Now compute
either MAP or
Bayesian estimate

$$P(B|\text{data}) = ?$$

$$+ \text{data} =$$

Parameter Estimation and Bayesian Networks



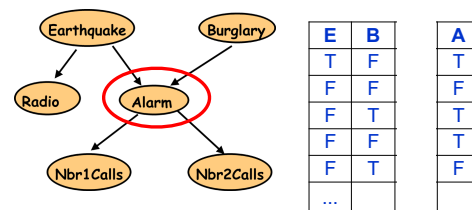
Prior

$$P(B|\text{data}) = \text{Beta}(1,4) \text{ "+ data" } = (3,7)$$

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

B	$\neg B$
.3	.7

Parameter Estimation and Bayesian Networks



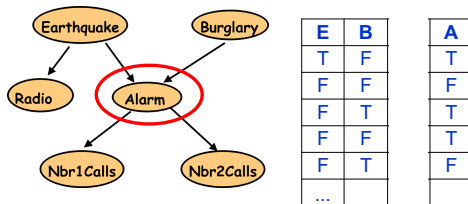
$$P(A|E,B) = ?$$

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

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

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

Parameter Estimation and Bayesian Networks



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

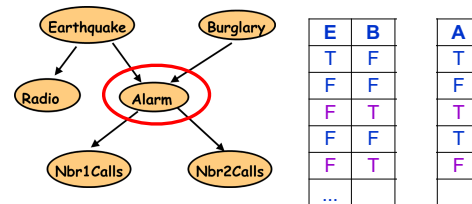
Prior

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

$$P(A|\neg E,B) = ? \text{ Beta}(2,3)$$

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

Parameter Estimation and Bayesian Networks



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

Prior

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

$$P(A|\neg E,B) = ? \text{ Beta}(2,3) \text{ + data} = (3,4)$$

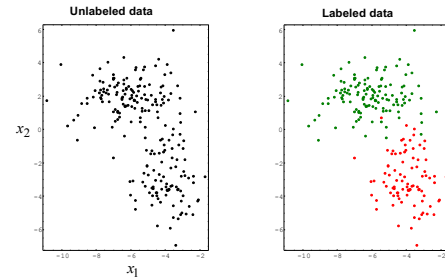
$$P(A|\neg E,\neg B) = ?$$

Expectation Maximization and Gaussian Mixtures

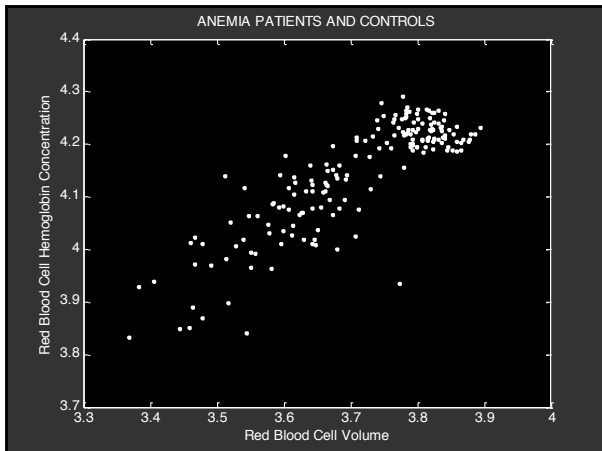
CSE 473

The problem of finding labels for unlabeled data

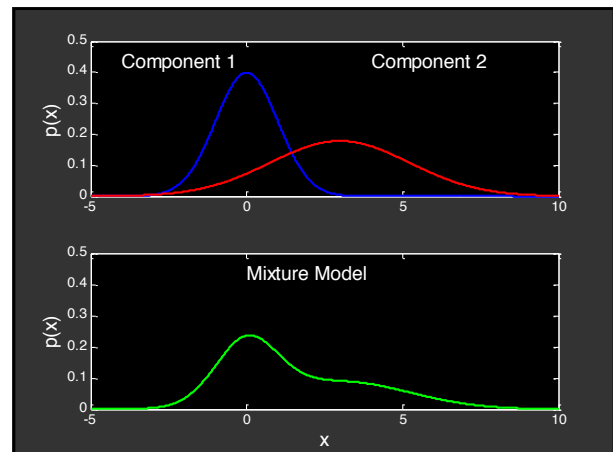
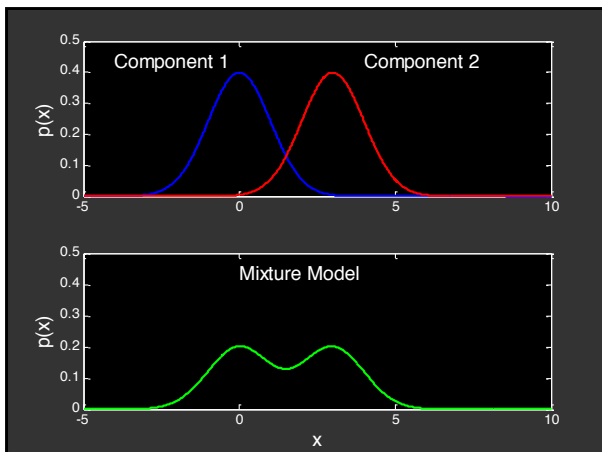
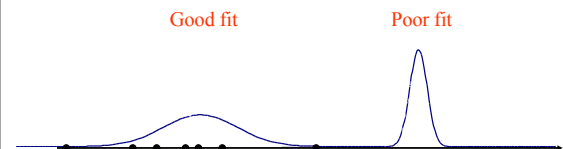
In nature, items often do not come with labels. How can we learn labels without a teacher?

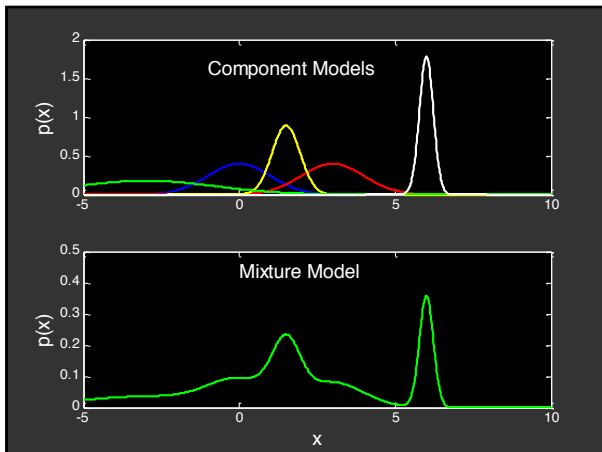


From Shadmehr & Diederichsen



Fitting a Gaussian PDF to Data





Mixtures

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

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

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

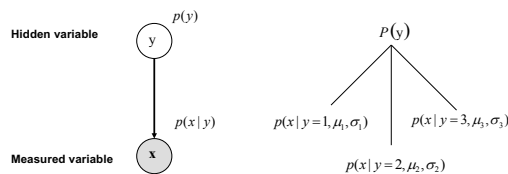
$$p(\mathbf{x}) = \sum_{j=1}^m P(y=j) p(\mathbf{x}|y=j)$$

We might hypothesize that the distributions are Gaussian:

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

$$p(\mathbf{x}|\theta) = \sum_{j=1}^m \underbrace{P(y=j)}_{\text{Mixing proportions}} \underbrace{N(\mathbf{x}|\mu_j, \Sigma_j)}_{\text{Normal distribution}}$$

Bayes Net for 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 $\Theta = \{\mu_1, \sigma_1, \mu_2, \sigma_2, \alpha_1\}$

Given data $D = \{x_1, \dots, x_N\}$, we want to find the parameters Θ 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:

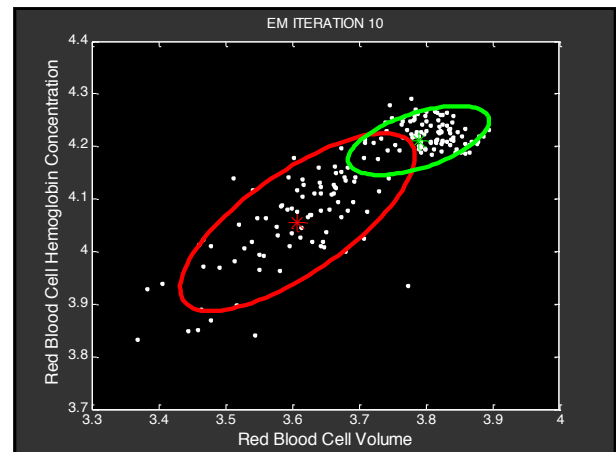
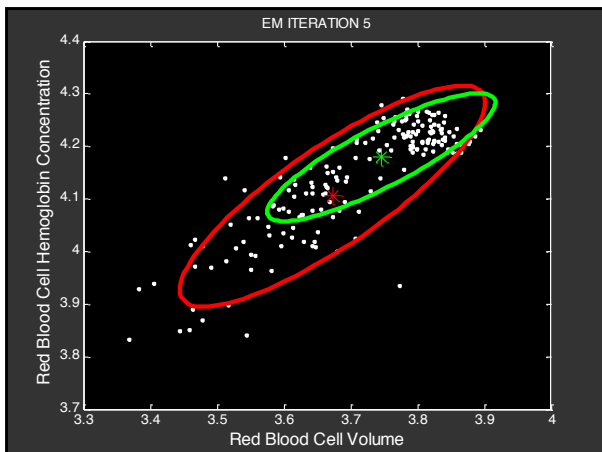
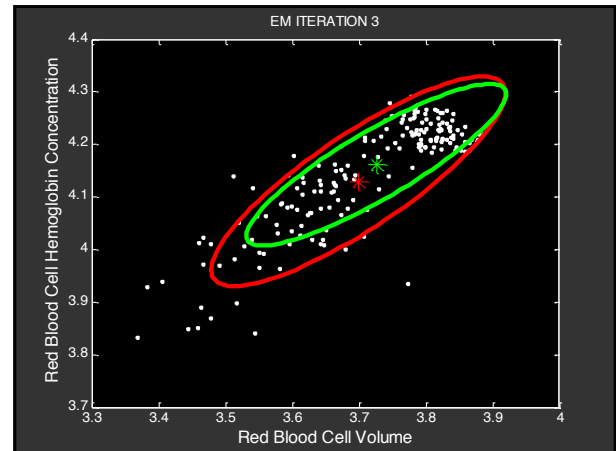
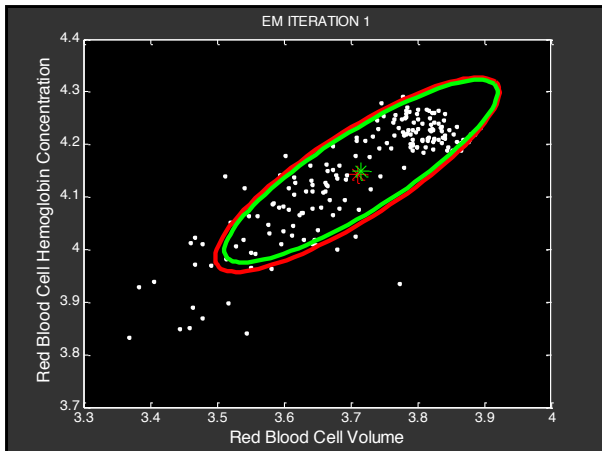
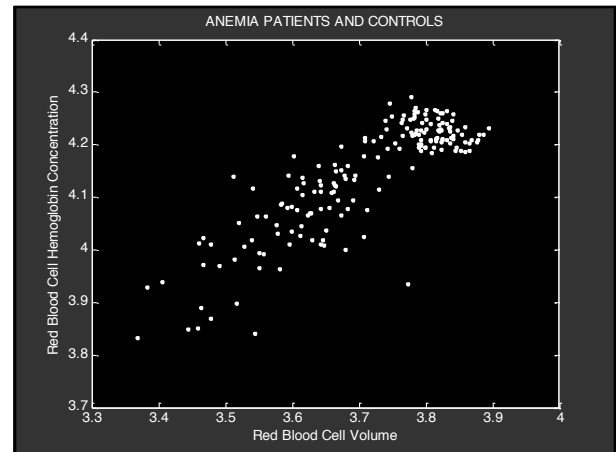
$$\mu_i \leftarrow \sum_j p_{ij} x_j / p_i$$

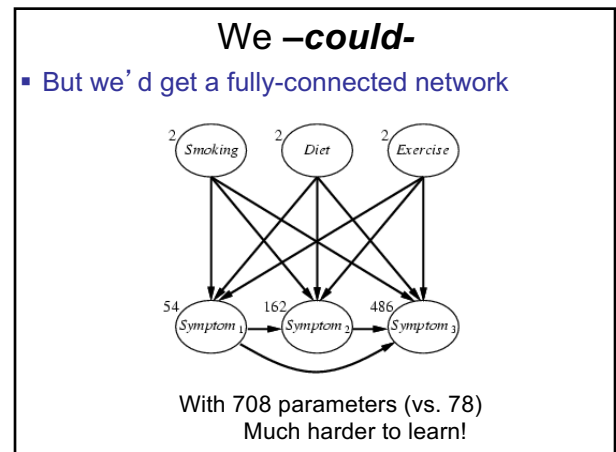
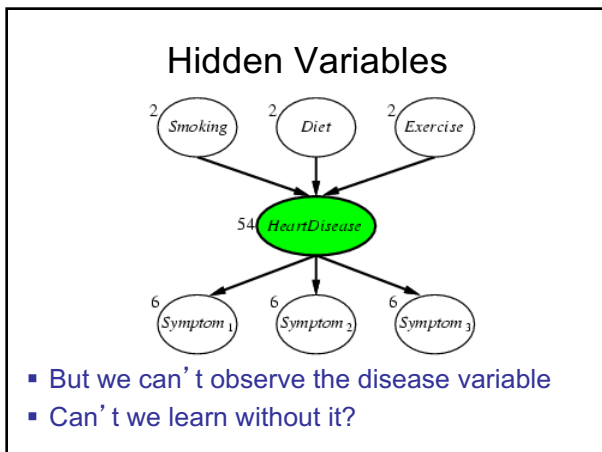
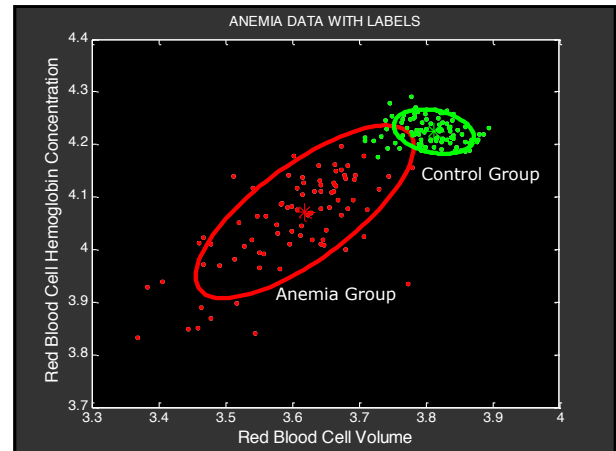
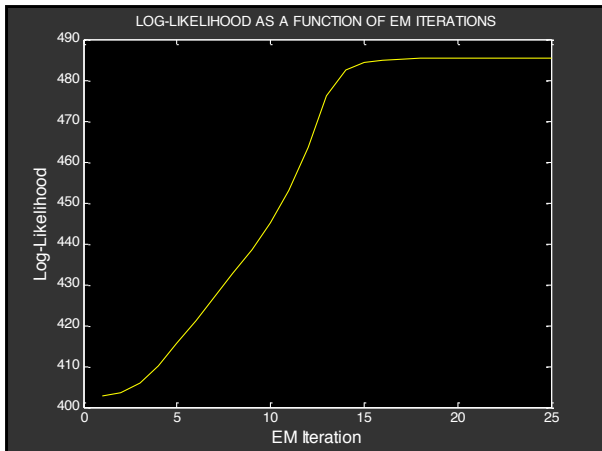
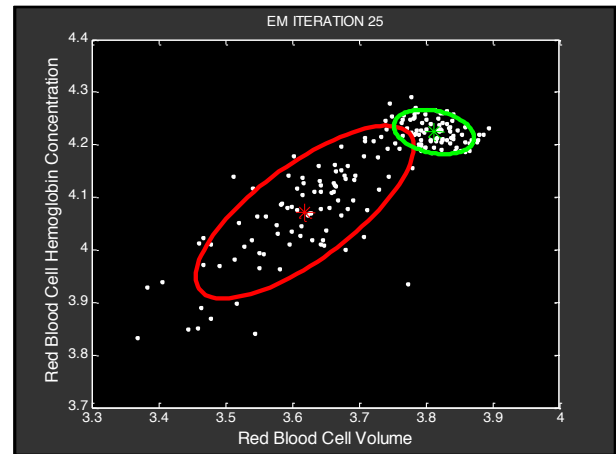
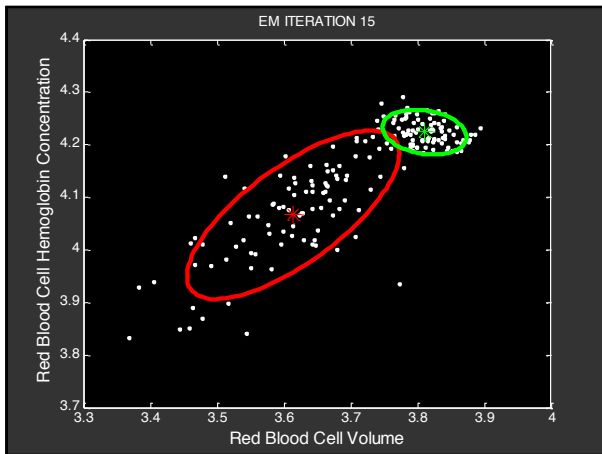
$$\sigma^2 \leftarrow \sum_j p_{ij} (x_j - \mu_i)^2 / p_i$$

$$w_i \leftarrow p_i$$

© D. Weld and D. Fox

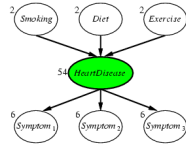
25





Chicken & Egg Problem

- If we knew that a training instance (patient) had the disease, then it'd be easy to learn $P(\text{symptom} \mid \text{disease})$
- If we knew params, e.g. $P(\text{symptom} \mid \text{disease})$ then it'd be easy to estimate if the patient had the disease



Expectation Maximization (EM) (high-level version)

- Pretend we **do** know the parameters
 - Initialize randomly
- **[E step]** Compute probability of instance having each possible value of the hidden variable
- **[M step]** Treating each instance as fractionally having both values compute the new parameter values
- Iterate until convergence!

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 for model complexity
Bayesian Information Criterion (BIC)

$P(D \mid \text{BN})$ – penalty

Penalty = $\frac{1}{2} (\# \text{ parameters}) \log (\# \text{ data points})$

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
 1. Start with some network structure
 2. Try to make a change (add or delete or reverse edge)
 3. 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?

