

CS 1675 Introduction to Machine Learning

Lecture 17

Learning complex distributions: Hidden variables and missing values

Milos Hauskrecht
milos@cs.pitt.edu
5329 Sennott Square

Probabilistic inferences

- BBN models compactly the full joint distribution by taking advantage of existing independences between variables
 - Simplifies the representation and learning of a model
 - Can be used for the different **inference tasks**
-

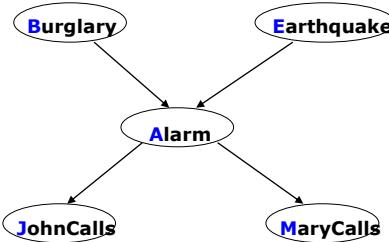
Inference in Bayesian network

- **Bad news:**

- Exact inference problem in BBNs is NP-hard (Cooper)
- Approximate inference is NP-hard (Dagum, Luby)

- **But** very often we can achieve significant improvements

- Assume our Alarm network



- Assume we want to compute: $P(J = T)$

Inference in Bayesian networks

Approach 2. Interleave sums and products

- Combines sums and product in a smart way (multiplications by constants can be taken out of the sum)

$$\begin{aligned}
 P(J = T) &= \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J = T | A = a) P(M = m | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e) \\
 &= \sum_{b \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J = T | A = a) P(M = m | A = a) P(B = b) \left[\sum_{e \in T, F} P(A = a | B = b, E = e) P(E = e) \right] \\
 &= \sum_{a \in T, F} P(J = T | A = a) \left[\sum_{m \in T, F} P(M = m | A = a) \right] \left[\sum_{b \in T, F} P(B = b) \left[\sum_{e \in T, F} P(A = a | B = b, E = e) P(E = e) \right] \right]
 \end{aligned}$$

Computational cost:

Number of additions: $1 + 2 * [1 + 1 + 2 * 1] = 9$

Number of products: $2 * [2 + 2 * (1 + 2 * 1)] = 16$

Variable elimination

- **Variable elimination:**

- Similar idea but interleave sum and products one variable at the time during inference
- E.g. Query $P(J = T)$ requires to eliminate A,B,E,M and this can be done in different order

$$\begin{aligned} P(J = T) &= \\ &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J = T | A = a) P(M = m | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e) \end{aligned}$$

Variable elimination

Assume order: M, E, B, A to calculate $P(J = T)$

$$= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J = T | A = a) P(M = m | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e)$$

Variable elimination

Assume order: M, E, B, A to calculate $P(J = T)$

$$\begin{aligned} &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J = T | A = a) P(M = m | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e) \\ &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} P(J = T | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e) \left[\sum_{m \in T, F} P(M = m | A = a) \right] \end{aligned}$$

Variable elimination

Assume order: M, E, B, A to calculate $P(J = T)$

$$\begin{aligned} &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J = T | A = a) P(M = m | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e) \\ &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} P(J = T | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e) \left[\sum_{m \in T, F} P(M = m | A = a) \right] \\ &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} P(J = T | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e) \boxed{1} \end{aligned}$$

Variable elimination

Assume order: M, E, B, A to calculate $P(J = T)$

$$\begin{aligned}
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J = T | A = a) P(M = m | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e) \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} P(J = T | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e) \left[\sum_{m \in T, F} P(M = m | A = a) \right] \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} P(J = T | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e) - 1 \quad \text{red arrow} \\
 &= \sum_{a \in T, F} \sum_{b \in T, F} P(J = T | A = a) P(B = b) \left[\sum_{e \in T, F} P(A = a | B = b, E = e) P(E = e) \right]
 \end{aligned}$$

Variable elimination

Assume order: M, E, B, A to calculate $P(J = T)$

$$\begin{aligned}
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J = T | A = a) P(M = m | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e) \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} P(J = T | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e) \left[\sum_{m \in T, F} P(M = m | A = a) \right] \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} P(J = T | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e) - 1 \quad \text{red arrow} \\
 &= \sum_{a \in T, F} \sum_{b \in T, F} P(J = T | A = a) P(B = b) \left[\sum_{e \in T, F} P(A = a | B = b, E = e) P(E = e) \right] \\
 &\qquad\qquad\qquad \boxed{\tau_1(A = a, B = b)} \quad \text{red box}
 \end{aligned}$$

$$\begin{aligned}
 \tau_1(A = a, B = b) = & \begin{array}{|c|c|} \hline A = T & A = F \\ \hline B = T & \begin{array}{|c|c|} \hline & \\ \hline & \\ \hline \end{array} \end{array} + \sum_{e \in T, F} P(A = F | B = T, E = e) P(E = e) \\
 & \begin{array}{|c|c|} \hline & \\ \hline & \\ \hline \end{array} + \sum_{e \in T, F} P(A = F | B = F, E = e) P(E = e)
 \end{aligned}$$

Variable elimination

Assume order: M, E, B, A to calculate $P(J=T)$

$$\begin{aligned}
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J=T | A=a) P(M=m | A=a) P(A=a | B=b, E=e) P(B=b) P(E=e) \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} P(J=T | A=a) P(A=a | B=b, E=e) P(B=b) P(E=e) \left[\sum_{m \in T, F} P(M=m | A=a) \right] \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} P(J=T | A=a) P(A=a | B=b, E=e) P(B=b) P(E=e) - 1 \quad \text{Red arrow points to } -1 \\
 &= \sum_{a \in T, F} \sum_{b \in T, F} P(J=T | A=a) P(B=b) \left[\sum_{e \in T, F} P(A=a | B=b, E=e) P(E=e) \right] \\
 &= \sum_{a \in T, F} \sum_{b \in T, F} P(J=T | A=a) P(B=b) \tau_1(A=a, B=b)
 \end{aligned}$$

Variable elimination

Assume order: M, E, B, A to calculate $P(J=T)$

$$\begin{aligned}
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J=T | A=a) P(M=m | A=a) P(A=a | B=b, E=e) P(B=b) P(E=e) \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} P(J=T | A=a) P(A=a | B=b, E=e) P(B=b) P(E=e) \left[\sum_{m \in T, F} P(M=m | A=a) \right] \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} P(J=T | A=a) P(A=a | B=b, E=e) P(B=b) P(E=e) - 1 \quad \text{Red arrow points to } -1 \\
 &= \sum_{a \in T, F} \sum_{b \in T, F} P(J=T | A=a) P(B=b) \left[\sum_{e \in T, F} P(A=a | B=b, E=e) P(E=e) \right] \\
 &= \sum_{a \in T, F} \sum_{b \in T, F} P(J=T | A=a) P(B=b) \tau_1(A=a, B=b) \\
 &= \sum_{a \in T, F} P(J=T | A=a) \left[\sum_{b \in T, F} P(B=b) \tau_1(A=a, B=b) \right]
 \end{aligned}$$

Variable elimination

Assume order: M, E, B, A to calculate $P(J=T)$

$$\begin{aligned}
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J=T | A=a) P(M=m | A=a) P(A=a | B=b, E=e) P(B=b) P(E=e) \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} P(J=T | A=a) P(A=a | B=b, E=e) P(B=b) P(E=e) \left[\sum_{m \in T, F} P(M=m | A=a) \right] \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} P(J=T | A=a) P(A=a | B=b, E=e) P(B=b) P(E=e) - 1 \quad \downarrow \\
 &= \sum_{a \in T, F} \sum_{b \in T, F} P(J=T | A=a) P(B=b) \left[\sum_{e \in T, F} P(A=a | B=b, E=e) P(E=e) \right] \quad \downarrow \\
 &= \sum_{a \in T, F} \sum_{b \in T, F} P(J=T | A=a) P(B=b) \tau_1(A=a, B=b) \\
 &= \sum_{a \in T, F} P(J=T | A=a) \left[\sum_{b \in T, F} P(B=b) \tau_1(A=a, B=b) \right] \\
 &\quad \downarrow \\
 &\boxed{\tau_2(A=a)}
 \end{aligned}$$

Variable elimination

Assume order: M, E, B, A to calculate $P(J=T)$

$$\begin{aligned}
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J=T | A=a) P(M=m | A=a) P(A=a | B=b, E=e) P(B=b) P(E=e) \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} P(J=T | A=a) P(A=a | B=b, E=e) P(B=b) P(E=e) \left[\sum_{m \in T, F} P(M=m | A=a) \right] \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} P(J=T | A=a) P(A=a | B=b, E=e) P(B=b) P(E=e) - 1 \quad \downarrow \\
 &= \sum_{a \in T, F} \sum_{b \in T, F} P(J=T | A=a) P(B=b) \left[\sum_{e \in T, F} P(A=a | B=b, E=e) P(E=e) \right] \quad \downarrow \\
 &= \sum_{a \in T, F} \sum_{b \in T, F} P(J=T | A=a) P(B=b) \tau_1(A=a, B=b) \\
 &= \sum_{a \in T, F} P(J=T | A=a) \left[\sum_{b \in T, F} P(B=b) \tau_1(A=a, B=b) \right] \quad \downarrow \\
 &= \sum_{a \in T, F} P(J=T | A=a) \quad \tau_2(A=a)
 \end{aligned}$$

Variable elimination

Assume order: $\mathbf{M}, \mathbf{E}, \mathbf{B}, \mathbf{A}$ to calculate $P(J = T)$

$$\begin{aligned}
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} \sum_{m \in T, F} P(J = T | A = a) P(M = m | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e) \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} P(J = T | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e) \left[\sum_{m \in T, F} P(M = m | A = a) \right] \\
 &= \sum_{b \in T, F} \sum_{e \in T, F} \sum_{a \in T, F} P(J = T | A = a) P(A = a | B = b, E = e) P(B = b) P(E = e) - 1 \quad \text{red arrow} \\
 &= \sum_{a \in T, F} \sum_{b \in T, F} P(J = T | A = a) P(B = b) \left[\sum_{e \in T, F} P(A = a | B = b, E = e) P(E = e) \right] \quad \text{red arrow} \\
 &= \sum_{a \in T, F} \sum_{b \in T, F} P(J = T | A = a) P(B = b) \tau_1(A = a, B = b) \\
 &= \sum_{a \in T, F} P(J = T | A = a) \left[\sum_{e \in T, F} P(B = b) \tau_1(A = a, B = b) \right] \quad \text{red arrow} \\
 &= \sum_{a \in T, F} P(J = T | A = a) \tau_2(A = a) = \boxed{P(J = T)}
 \end{aligned}$$

Inference in Bayesian network

- **Exact inference algorithms:**
 - Variable elimination
 - Recursive decomposition (Cooper, Darwiche)
 - Symbolic inference (D'Ambrosio)
 - Belief propagation algorithm (Pearl)
 - Clustering and joint tree approach (Lauritzen, Spiegelhalter)
 - Arc reversal (Olmsted, Schachter)

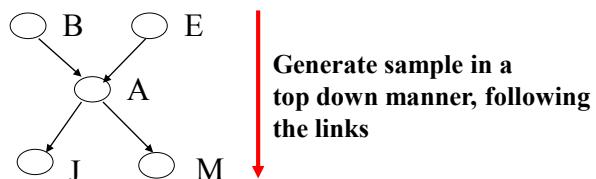
- **Approximate inference algorithms:**
 - Monte Carlo methods:
 - Forward sampling, Likelihood sampling
 - Variational methods

Monte Carlo approaches

- **MC approximation:**
 - The probability is approximated using **sample frequencies**
 - **Example:**

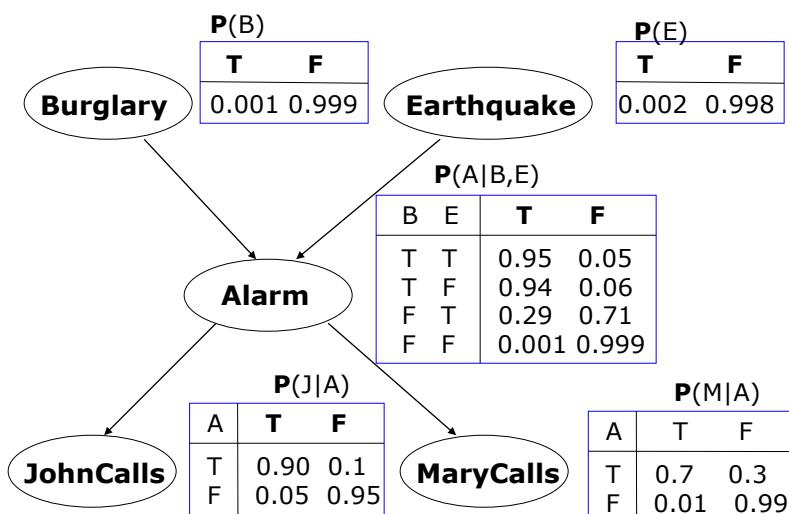
$$\tilde{P}(B = T, J = T) = \frac{N_{B=T, J=T}}{N}$$

#samples with $B = T, J = T$
total # samples
- **Sample generation:** BBN sampling of the joint is easy

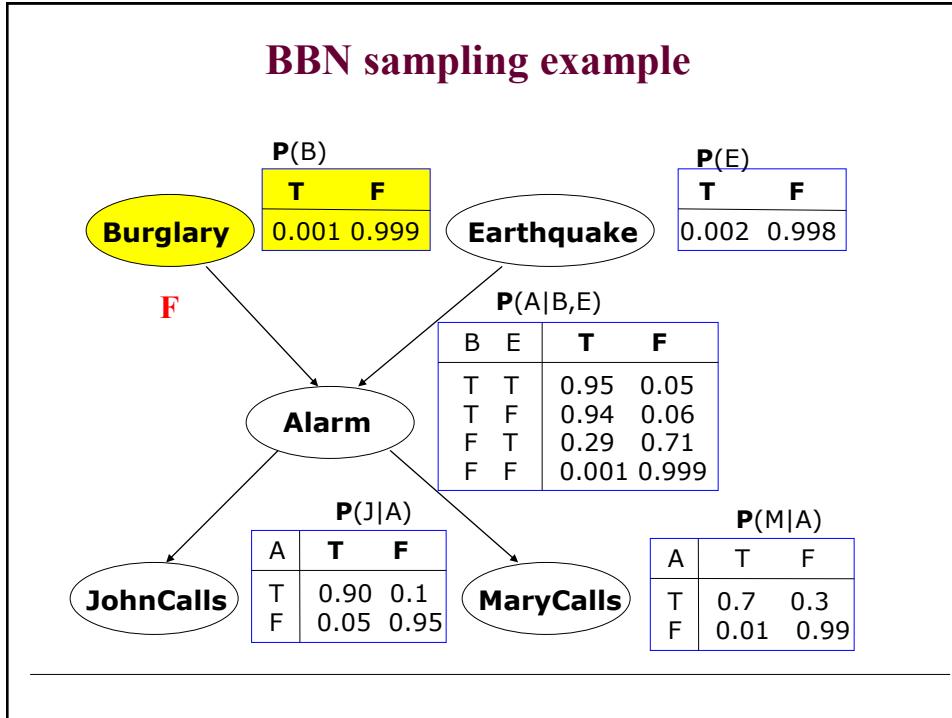


- One sample gives one assignment of values to all variables

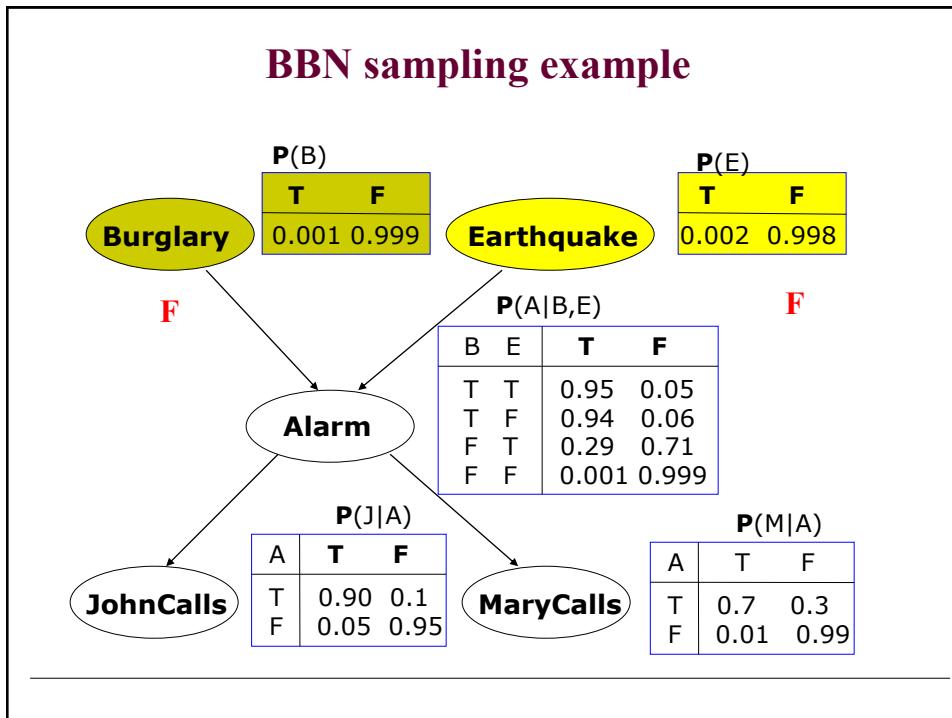
BBN sampling example



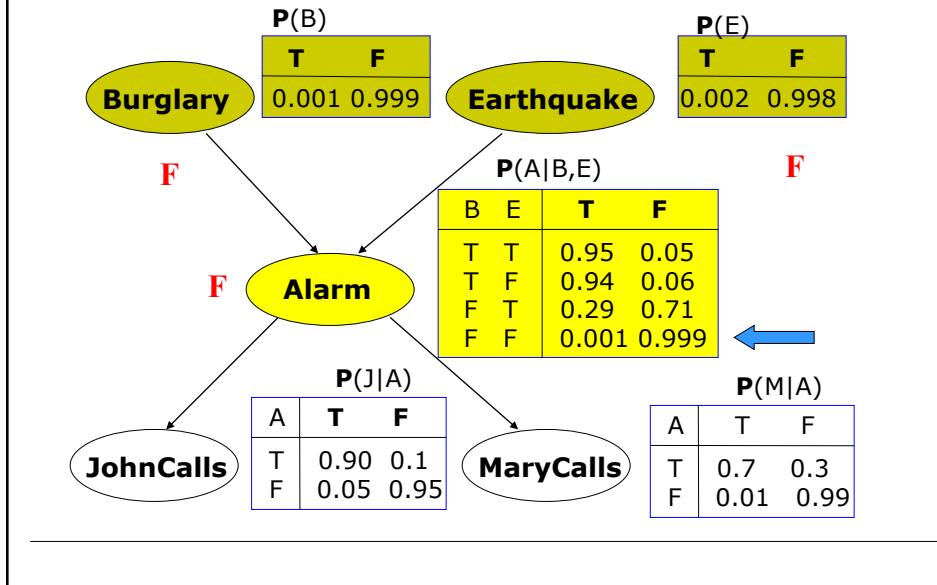
BBN sampling example



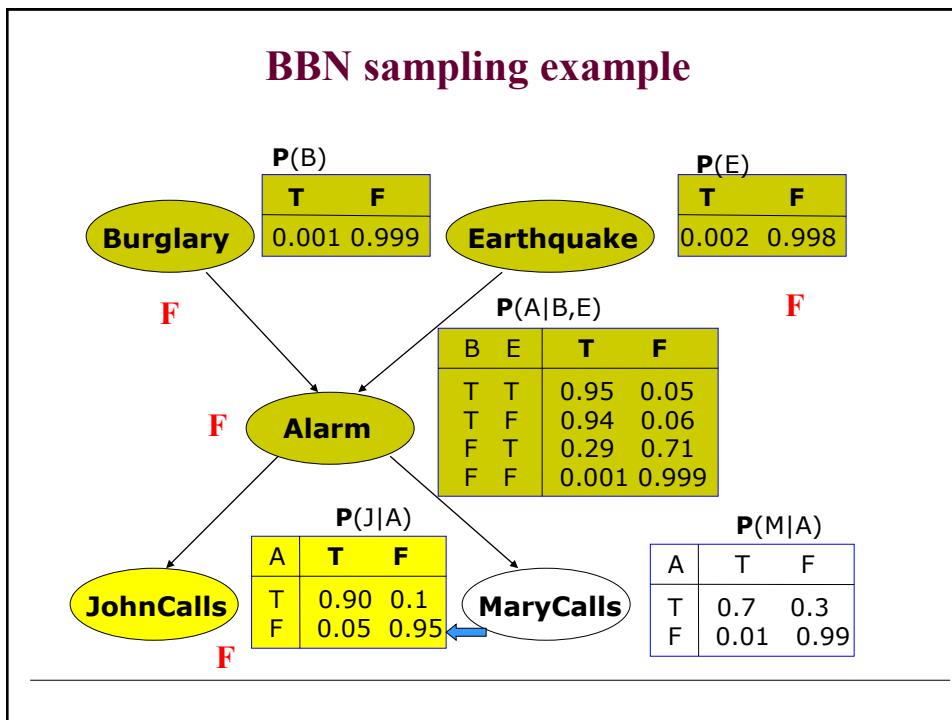
BBN sampling example



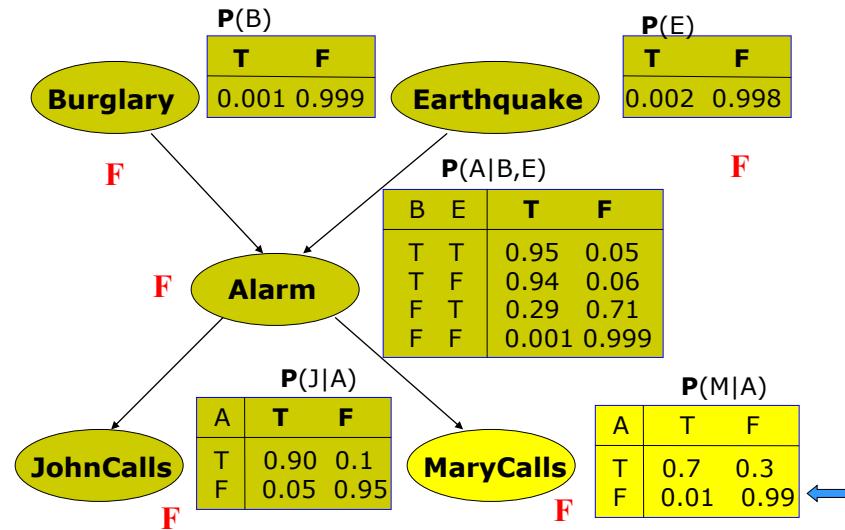
BBN sampling example



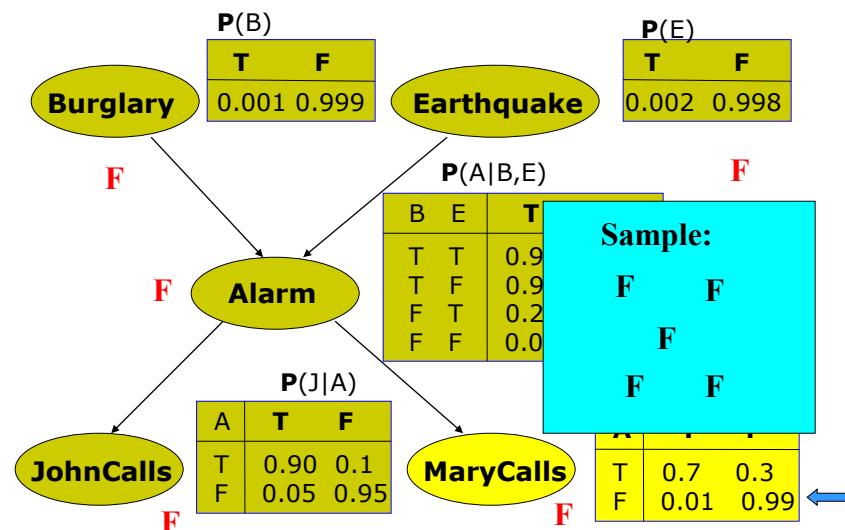
BBN sampling example



BBN sampling example



BBN sampling example



Monte Carlo approaches

- **MC approximation of conditional probabilities:**

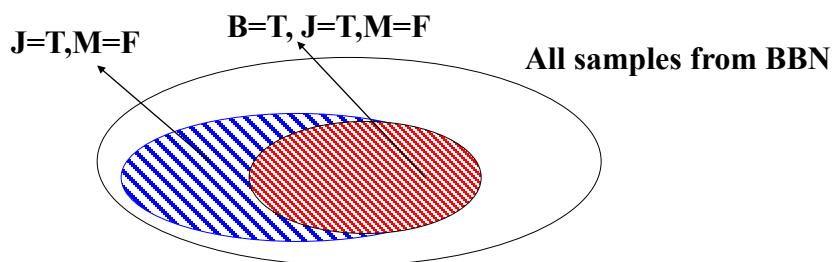
– The probability is approximated using sample frequencies

– **Example:**

samples with $B = T, J = T, M = F$

$$\tilde{P}(B = T \mid J = T, M = F) = \frac{N_{B=T, J=T, M=F}}{N_{J=T, M=F}}$$

samples with $J = T, M = F$

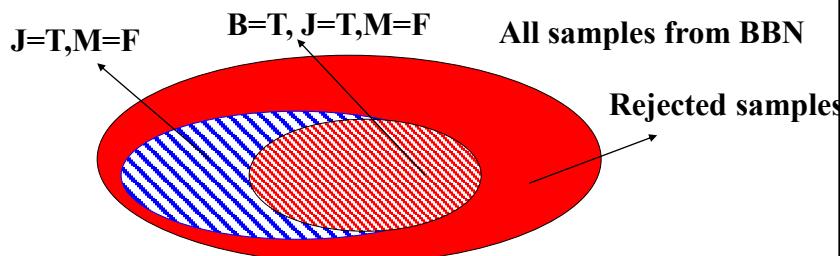


Monte Carlo approaches

- **Rejection sampling**

- Generate samples from the full joint by sampling BBN
- Use only samples that agree with the condition, the remaining samples are rejected

- **Problem:** many samples can be rejected



Likelihood weighting

Idea: generate only samples consistent with an evidence (or conditioning event)

- Benefit: Avoids inefficiencies of rejection sampling

Problem:

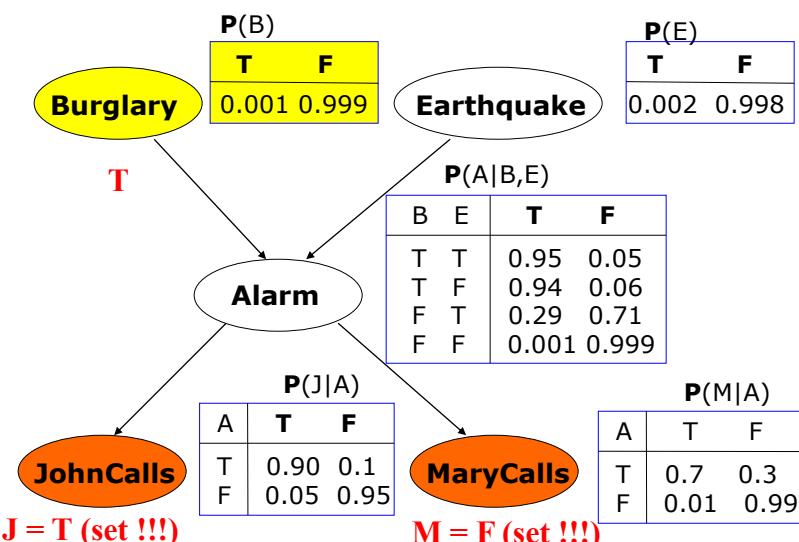
- the distribution generated by enforcing the conditioning variables to set values is biased
- simple counts are not sufficient to estimate the probabilities

Solution:

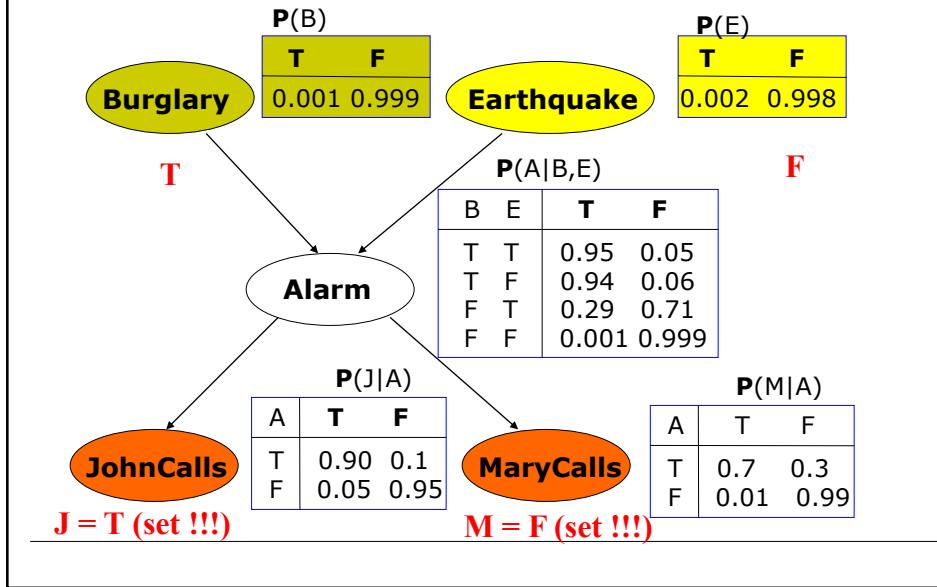
- With every sample keep a weight with which it should count towards the estimate

$$\tilde{P}(B = T \mid J = T, M = F) = \frac{\sum_{\text{samples with } B=T, M=F \text{ and } J=T} w_{B=T|J=T, M=F}}{\sum_{\text{samples with any value of } B \text{ and } J=T, M=F} w_{B=x|J=T, M=F}}$$

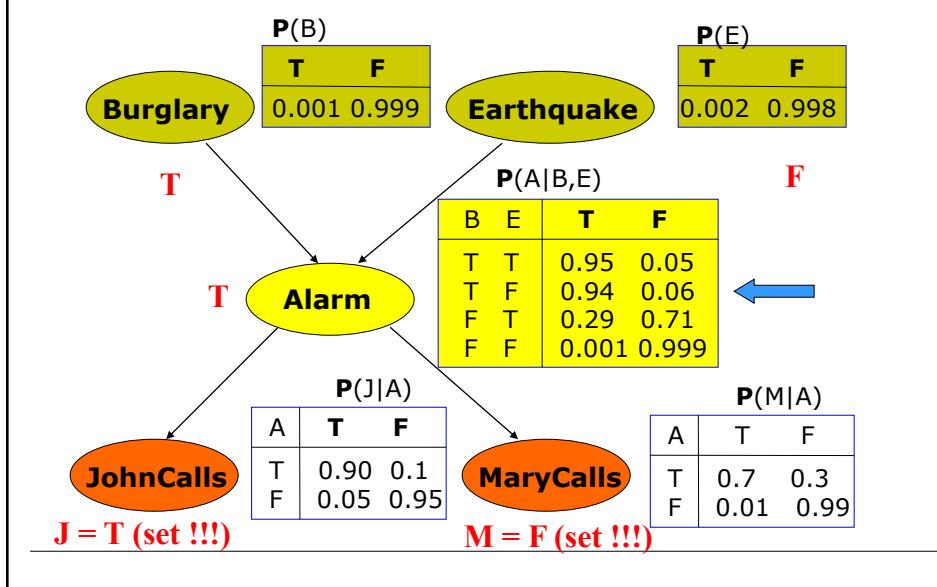
BBN likelihood weighting example



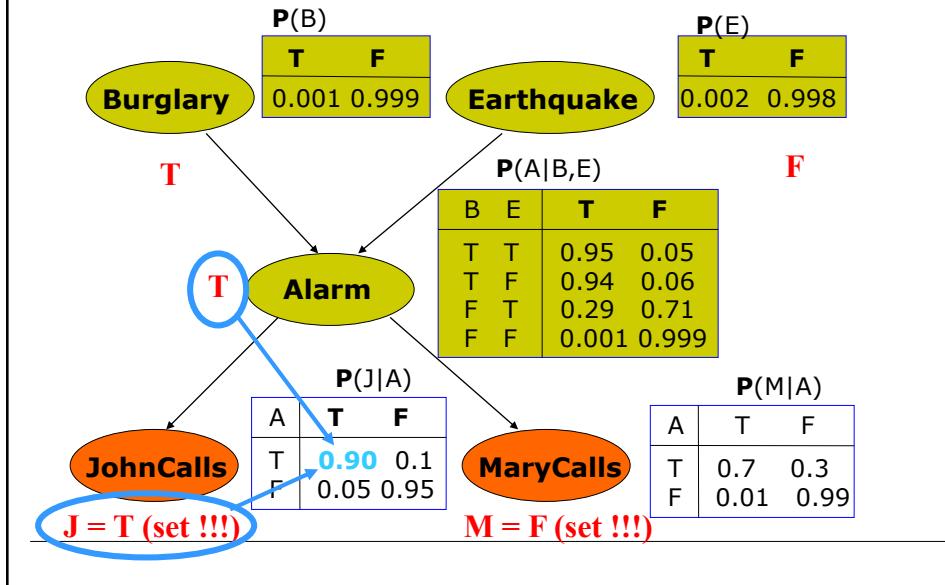
BBN likelihood weighting example



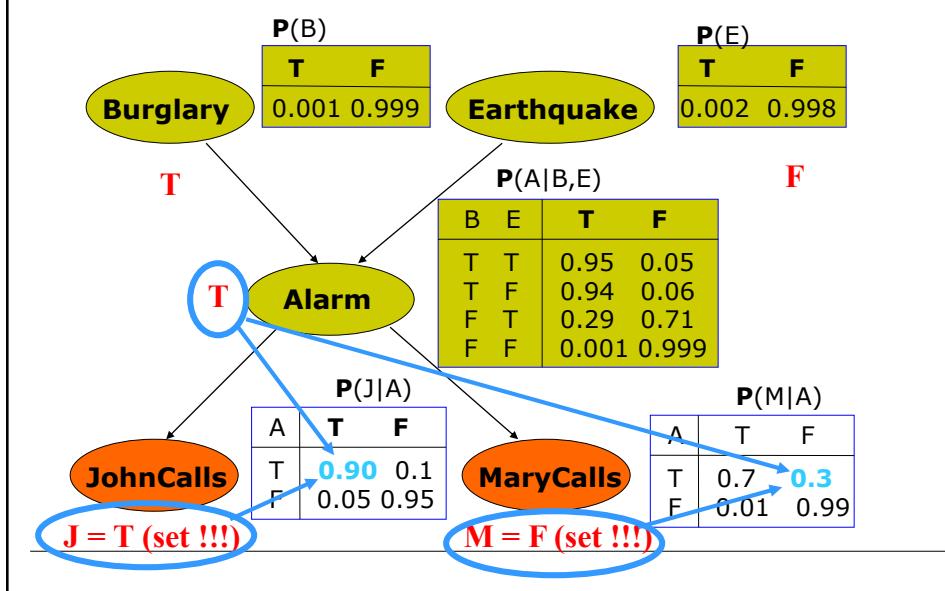
BBN likelihood weighting example



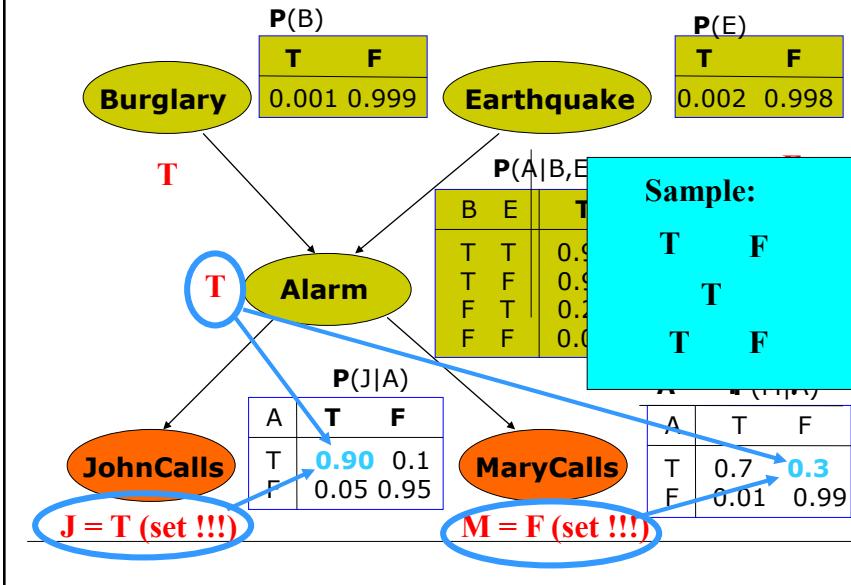
BBN likelihood weighting example



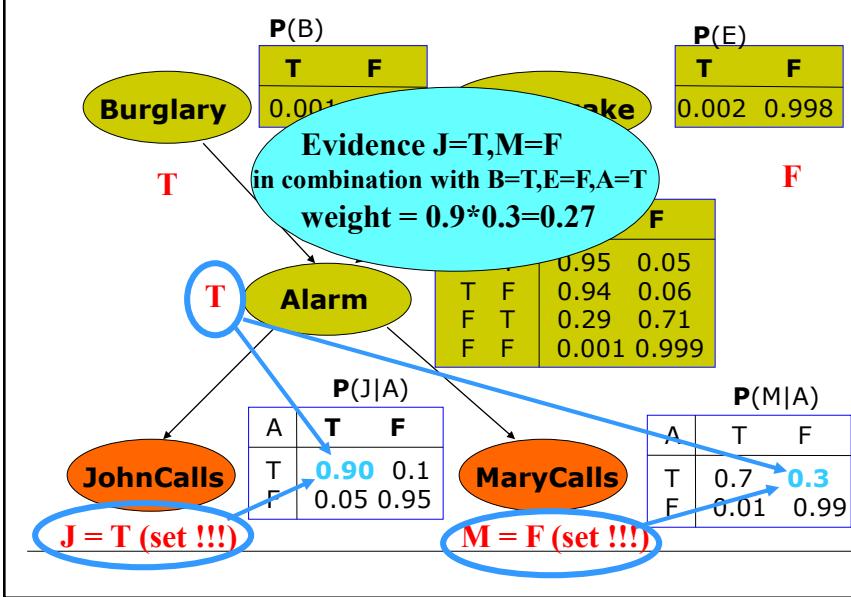
BBN likelihood weighting example



BBN likelihood weighting example

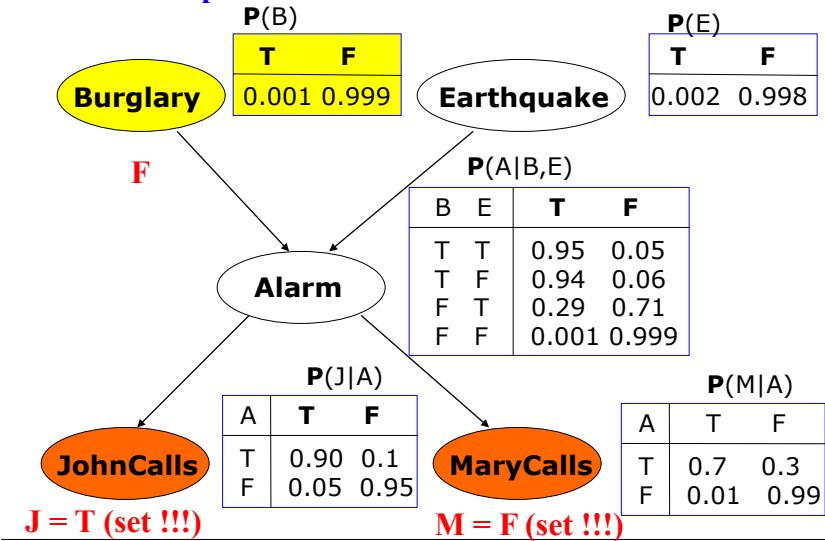


BBN likelihood weighting example



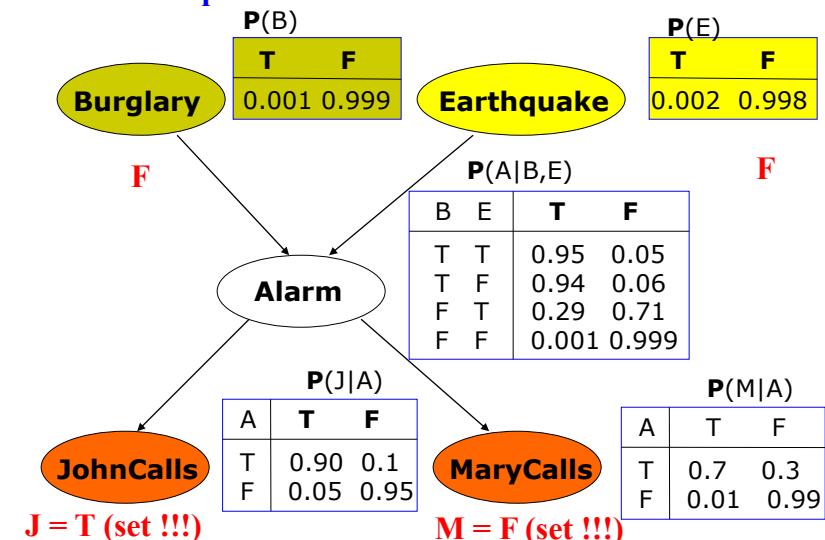
BBN likelihood weighting example

Second sample



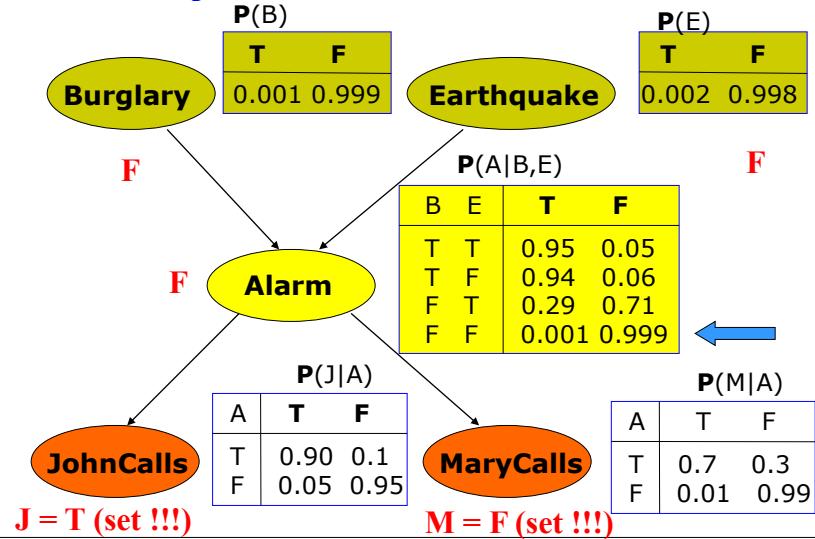
BBN likelihood weighting example

Second sample



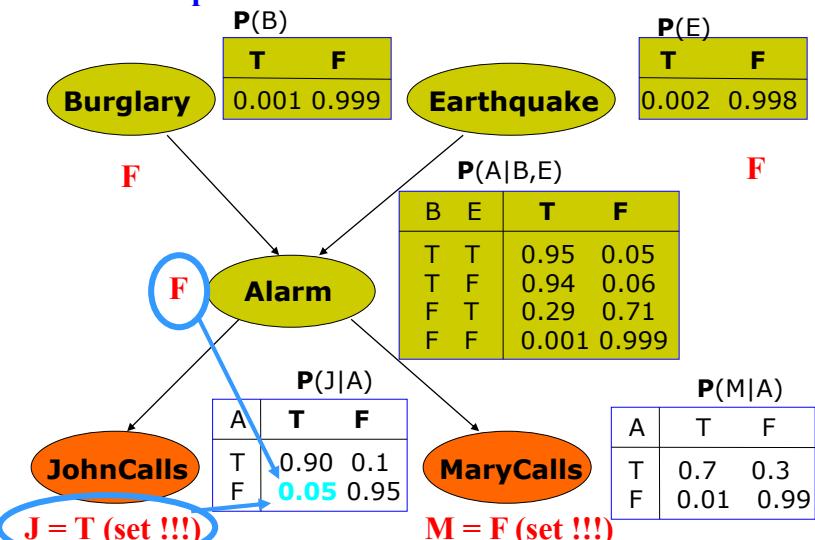
BBN likelihood weighting example

Second sample



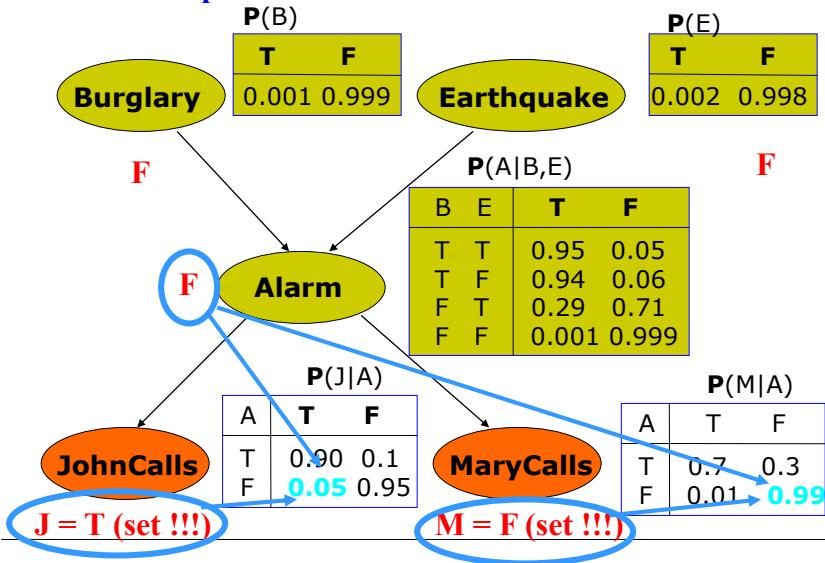
BBN likelihood weighting example

Second sample



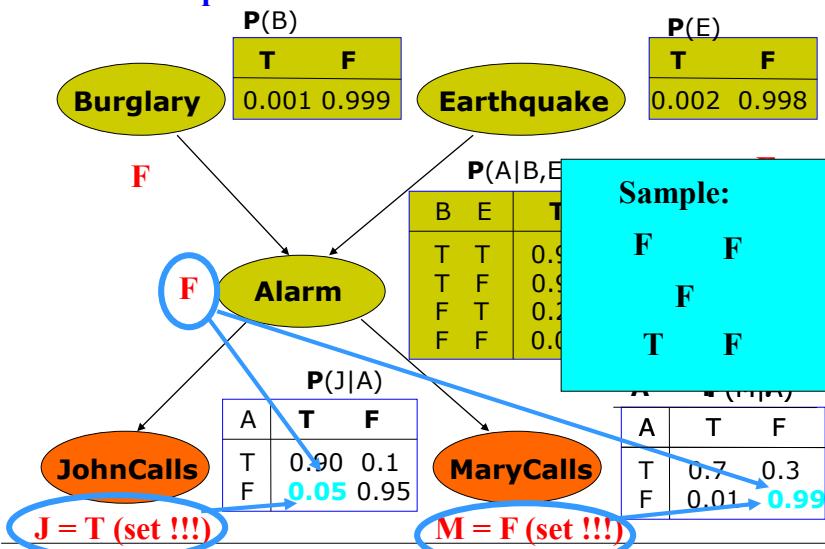
BBN likelihood weighting example

Second sample



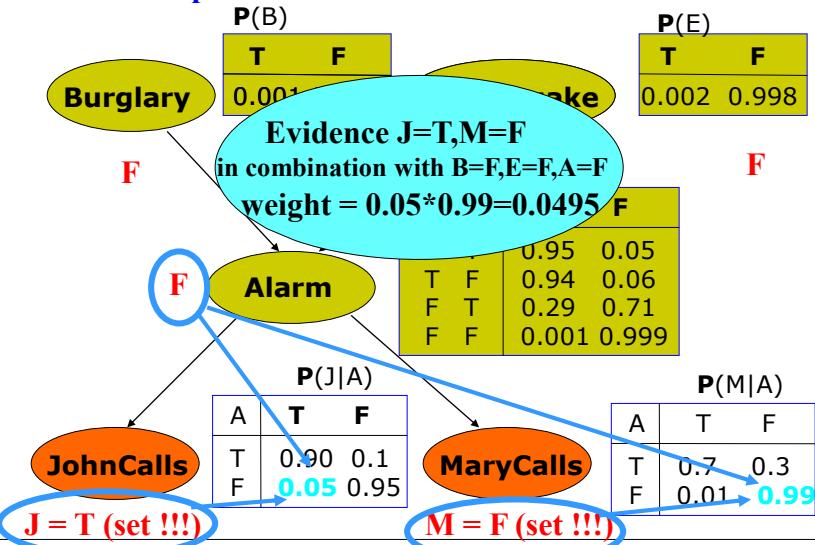
BBN likelihood weighting example

Second sample



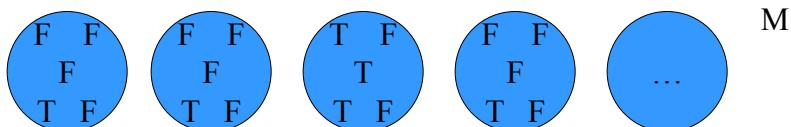
BBN likelihood weighting example

Second sample



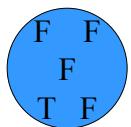
Likelihood weighting

- Assume we have generated the following M samples:



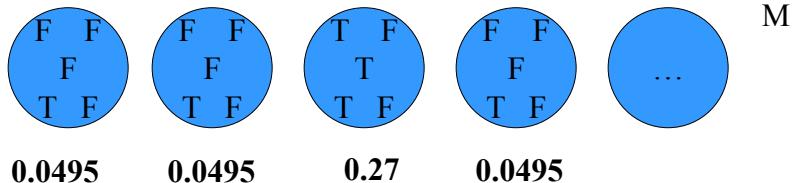
How to make the samples consistent?

Weight each sample by probability with which it agrees with the conditioning evidence $P(e)$.



Likelihood weighting

- Assume we have generated the following M samples:



$$\tilde{P}(B = T \mid J = T, M = F) = \frac{\sum_{\substack{\text{samples with } B=T, M=F \text{ and } J=T}} w_{B=T|J=T, M=F}}{\sum_{\substack{\text{samples with any value of } B \text{ and } J=T, M=F}} w_{B=x|J=T, M=F}}$$

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 X with parameters Θ
- 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

Hidden variables

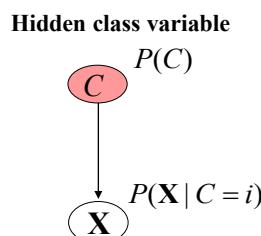
Modeling assumption:

Observed Variables $\mathbf{X} = \{X_1, X_2, \dots, X_n\}$

- Additional **(hidden) variables** may be added to the model
 - they are never observed in data

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

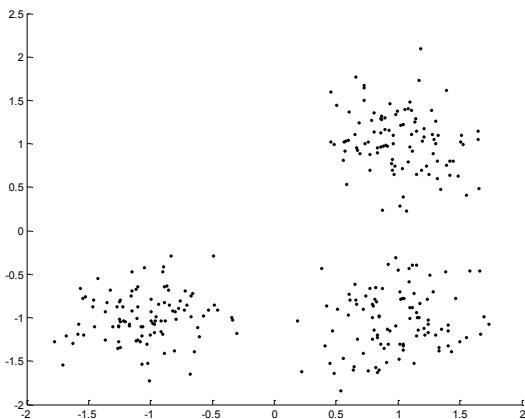


Example:

- Latent variable models
 - hidden classes (categories)

Gaussian mixture model

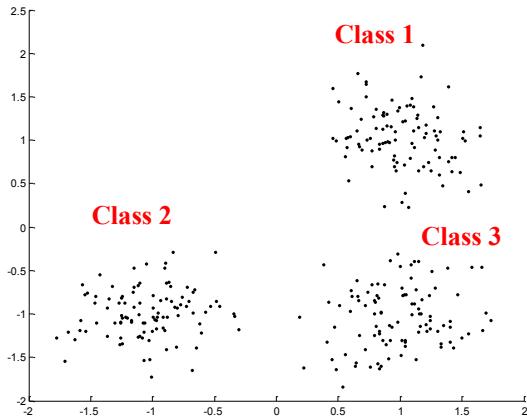
Data: 2 dimensional space of \mathbf{X}



Gaussian mixture model

Data: 2 dimensional space of \mathbf{X}

As if data came from 3 different classes – each with its own density

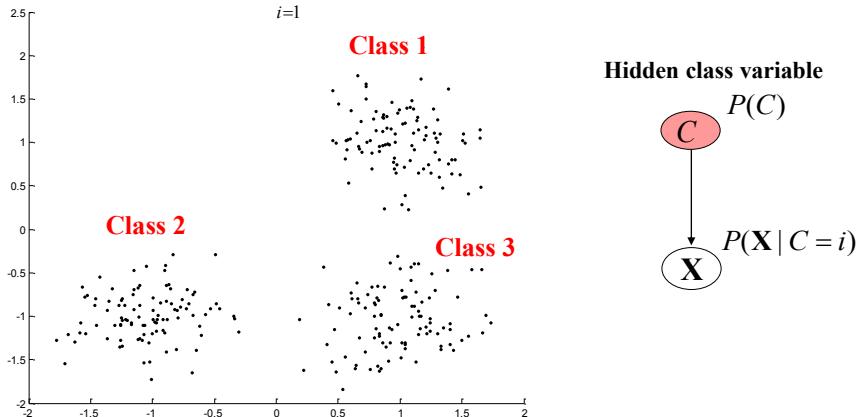


Gaussian mixture model

Assumption: data are coming from multiple Gaussians

- **Hidden variable:** models the different classes

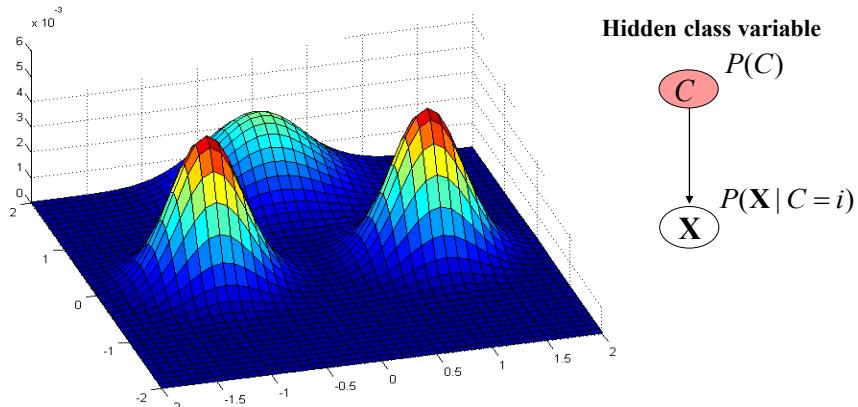
$$P(\mathbf{X}) = \sum_{i=1}^3 P(C=i)P(\mathbf{X} | C=i)$$



Mixture of Gaussians

- Density function for the Mixture of Gaussians model:

$$P(\mathbf{X}) = \sum_{i=1}^3 P(C=i)P(\mathbf{X} | C=i)$$



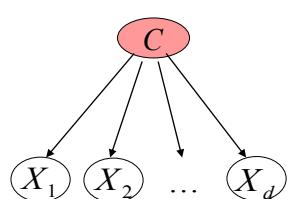
Naïve Bayes with a hidden class variable

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

Example:

- Naïve Bayes model with a hidden class variable

Hidden class variable



Attributes are independent given the class

$$P(\mathbf{X}) = \sum_{i=1}^k P(C=i)P(\mathbf{X} | C=i)$$

$$P(\mathbf{X}) = \sum_{i=1}^k P(C=i) \prod_{j=1}^d P(X_j | C=i)$$

- Useful in customer profiles**

– Class value = type of customers

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

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

Missing values of x_1^{i+1}, x_2^{i+1}

Etc.

- **Example: medical records**

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

Density estimation

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

ML estimate:

$$\max_{\Theta} p(D | \Theta, \xi)$$

Can be optimized but it is hard with hidden vars and missing values

- **Optimization methods for ML:** gradient-ascent, conjugate gradient, Newton-Raphson, etc.

Problem: No or very small advantage from the structure of the corresponding belief network when there are unobserved values

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**

General EM

Iterative re-estimation of parameters

Pick initial set of model parameters Θ

Repeat

Set $\Theta' = \Theta$

Expectation step. For all hidden and missing variables (and their possible value assignments) calculate their expectations for all data instances given the parameters Θ'

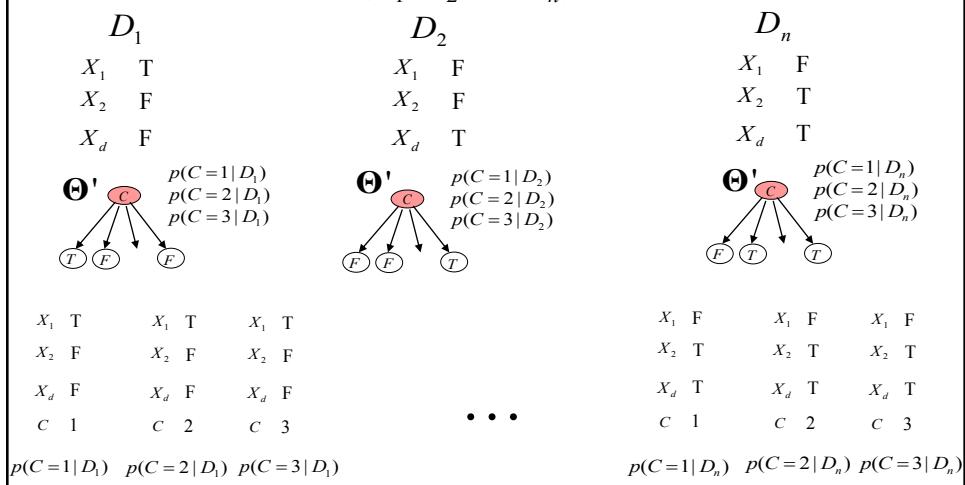
Maximization step. Compute the new estimates of Θ by considering the expectations of the different value completions (for hidden variables and missing values for all instances)

Till no improvement possible

EM example: expectation step

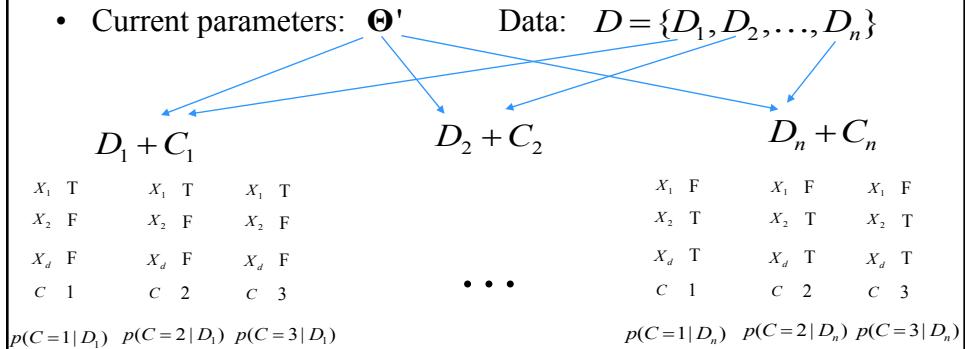
- Assume current set of parameters Θ'

$$D = \{D_1, D_2, \dots, D_n\}$$



EM example: maximization step

- Current parameters: Θ'



New parameters (one step of EM iteration)

$$\begin{aligned}\Theta &= \arg \max_{\Theta} E_{P(C|D, \Theta')} \log P(D, C | \Theta) \\ &= \arg \max_{\Theta} E_{P(C|D, \Theta')} \log \prod_{i=1}^n P(D_i, C_i | \Theta)\end{aligned}$$

EM advantages

Key advantages:

- In many problems (e.g. Bayesian belief networks)
 - the maximization step can be carried out in the closed form
 - We directly optimize using quantities corresponding to expected counts
- Climbs the gradient, but it does not need a learning rate, automatically renormalized update

Example: Gaussian mixture model

Probability of 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

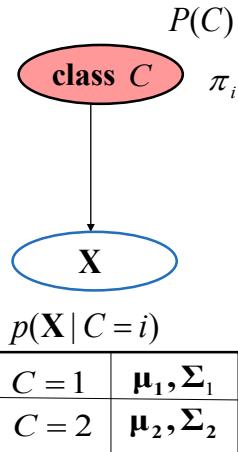
$$p(C=i)$$

= 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 a Gaussian for class i)



Special feature: C is hidden !!!

Example: Gaussian mixture model

Assume a generative classifier model based on the QDA:

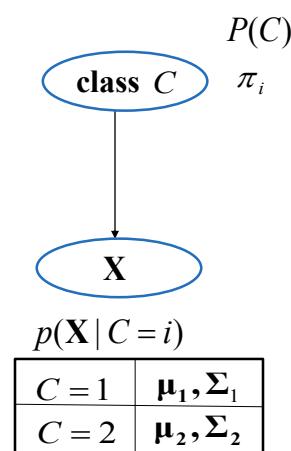
- **The class labels are known.** The ML estimate is

$$N_i = \sum_{j:C_j=i} 1$$

$$\tilde{\pi}_i = \frac{N_i}{N}$$

$$\tilde{\boldsymbol{\mu}}_i = \frac{1}{N_i} \sum_{j:C_j=i} \mathbf{x}_j$$

$$\tilde{\boldsymbol{\Sigma}}_i = \frac{1}{N_i} \sum_{j:C_j=i} (\mathbf{x}_j - \tilde{\boldsymbol{\mu}}_i)(\mathbf{x}_j - \tilde{\boldsymbol{\mu}}_i)^T$$



Example: Gaussian mixture model

- In the Gaussian mixture **Gaussians are not labeled**
- We can apply **EM algorithm**:

– re-estimation based on the class posterior

$$h_{il} = p(C_l = i | \mathbf{x}_l, \Theta') = \frac{p(C_l = i | \Theta') p(x_l | C_l = i, \Theta')}{\sum_{u=1}^m p(C_l = u | \Theta') p(x_l | C_l = u, \Theta')}$$

$$N_i = \sum_l h_{il}$$

$$\tilde{\pi}_i = \frac{N_i}{N}$$

$$\tilde{\mu}_i = \frac{1}{N_i} \sum_l h_{il} \mathbf{x}_j$$

$$\tilde{\Sigma}_i = \frac{1}{N_i} \sum_l h_{il} (\mathbf{x}_j - \tilde{\mu}_i)(\mathbf{x}_j - \tilde{\mu}_i)^T$$

Count replaced with the expected count