

## CS 3750 Machine Learning Lecture 4

### Monte Carlo methods

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

---

CS 3750 Advanced Machine Learning

### Monte Carlo inference

- Let us assume we have a probability distribution  $P(X)$  represented e.g. using BBN or MRF, and we want calculate  $P(X=x)$  ( $P(x)$  in short)
- We can use exact probabilistic inference, but it may be hard to calculate
- **Monte Carlo approximation:**
  - **Idea:** The probability  $P(x)$  is approximated using sample frequencies
- **Idea (first method):**
  - Generate a random sample  $D$  of size  $M$  from  $P(X)$
  - Estimate  $P(x)$  as:

$$\hat{P}_D(X = x) = \frac{M_{x=x}}{M}$$

---

CS 3750 Advanced Machine Learning

## Absolute Error Bound

- **Hoeffding's bound** lets us bound the probability with which the estimate  $\hat{P}_D(x)$  differs from  $P(x)$  by more than  $\varepsilon$

$$P(\hat{P}_D(x) \notin [P(x) - \varepsilon, P(x) + \varepsilon]) \leq 2e^{-2M\varepsilon^2} \leq \delta$$

The bound can be used to decide on how many samples are required to achieve a desired accuracy:

$$M \geq \frac{\ln(2/\delta)}{2\varepsilon^2}$$

---

3

## Relative Error Bound

- **Chernoff's bound** lets us bound the probability of the estimate  $\hat{P}_D(x)$  exceeding a relative error  $\varepsilon$  of the true value  $P(x)$

$$P(\hat{P}_D(x) \notin P(x)(1 \pm \varepsilon)) \leq 2e^{-MP(x)\varepsilon^2/3} \leq \delta$$

- This leads to the following sample complexity bound:

$$M \geq 3 \frac{\ln(2/\delta)}{P(x)\varepsilon^2}$$

---

4

## Monte Carlo inference challenges

### Challenge 1: How to generate $M$ (unbiased) examples from the target distribution $P(X)$ ?

- Generating (unbiased) examples from  $P(X)$  may be hard, or very inefficient

#### Example:

- Assume I have a distribution over 100 binary variables
  - There are  $2^{100}$  possible configurations of variable values
- **Trivial sampling solution:**
  - calculate and store the probability of each configuration
  - Pick randomly a configuration based on its probability
- **Problem:** terribly inefficient in time and memory

CS 3750 Advanced Machine Learning

## Monte Carlo inference challenges

### Challenge 2: How to estimate the expected value of $f(x)$ for $P(x)$ :

- Generally, we can estimate this expectation by generating samples  $x[1], \dots, x[M]$  from  $P$ , and then estimating it as:

$$E_P[f] = \sum_x P(x) f(x) \qquad E_P[f] = \int_x p(x) f(x) dx$$
$$\hat{\Phi} = \hat{E}_P[f] = \frac{1}{M} \sum_{m=1}^M f(x[m])$$

- Using the central limit theorem, the estimate  $\hat{\Phi}$  follows  $N\left(0, \frac{\sigma^2}{M}\right)$ 
  - Where is the variance for  $f(x)$  is

$$\sigma^2 = \int_x p(x) [f(x) - E_P(f(x))]^2 dx$$

- **Problem:** we are unable to efficiently sample  $P(x)$ . What to do?

CS 3750 Advanced Machine Learning

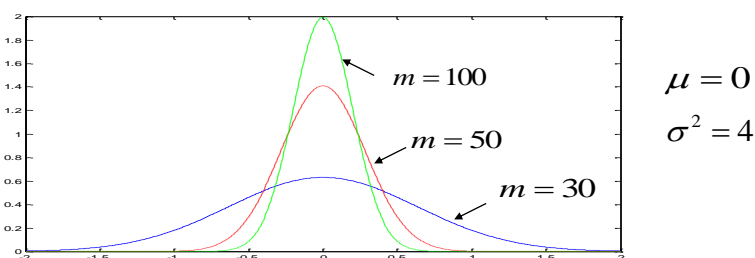
## Central limit theorem

- **Central limit theorem:**

Let random variables  $X_1, X_2, \dots, X_m$  form a random sample from a distribution with mean  $\mu$  and variance  $\sigma^2$ , then if the sample  $n$  is large, the distribution

$$\sum_{i=1}^m X_i \approx N(m\mu, m\sigma^2) \quad \text{or} \quad \frac{1}{m} \sum_{i=1}^m X_i \approx N(\mu, \sigma^2 / m)$$

Effect of increasing the sample size  $m$  on the sample mean:



CS 3750 Advanced Machine Learning

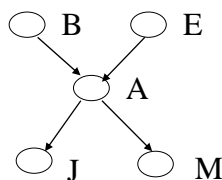
## Monte Carlo inference: BBNs

**Challenge 1: How to generate  $M$  (unbiased) examples from the target distribution  $P(X)$  defined by a BBN?**

- **Good news: Sample generation for the full joint defined by the BBN is easy**

- One top down sweep through the network lets us generate one example according to  $P(X)$

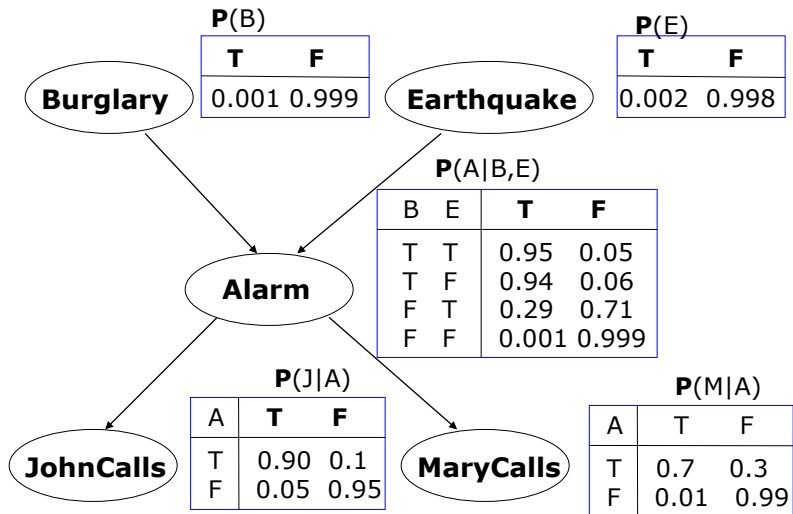
- **Example:**



Examples are generated in a top down manner, following the links

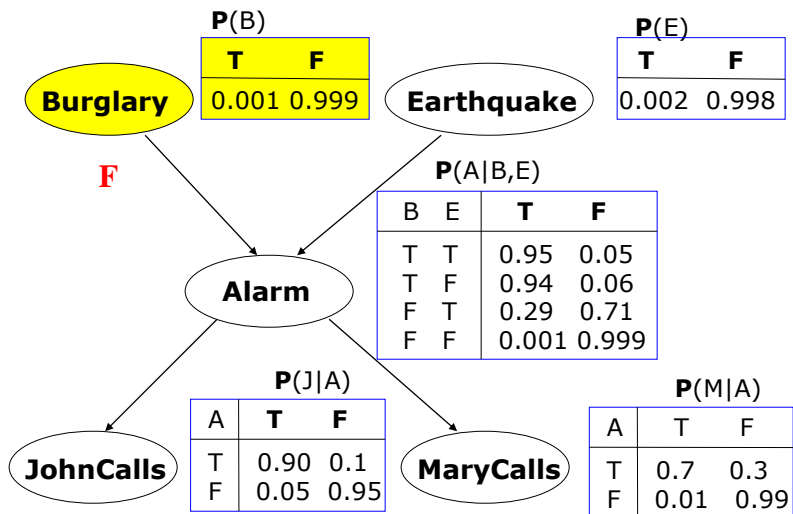
CS 3750 Advanced Machine Learning

## BBN sampling example



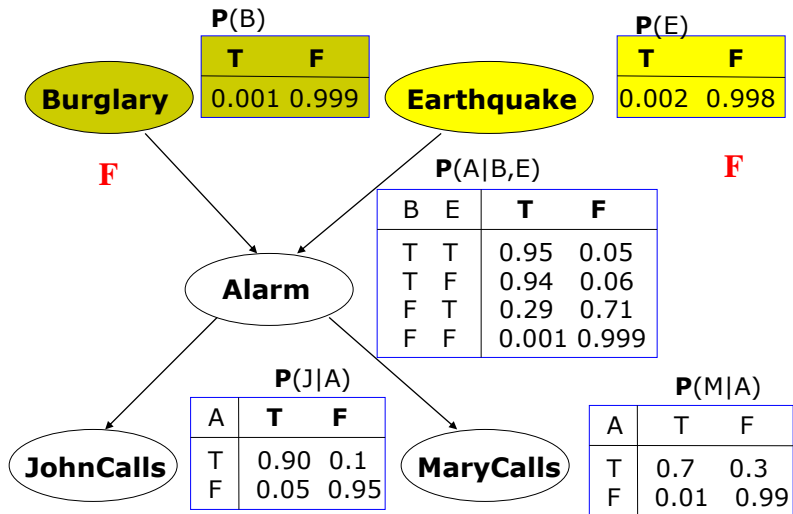
CS 3750 Advanced Machine Learning

## BBN sampling example



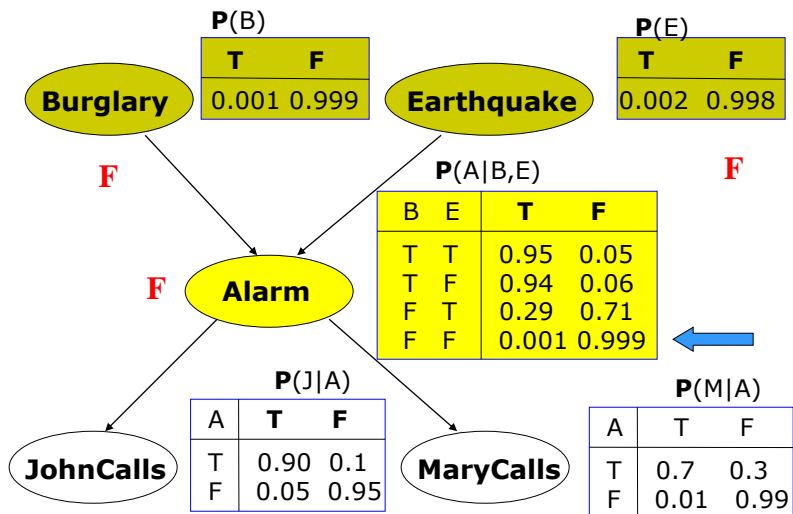
CS 3750 Advanced Machine Learning

## BBN sampling example



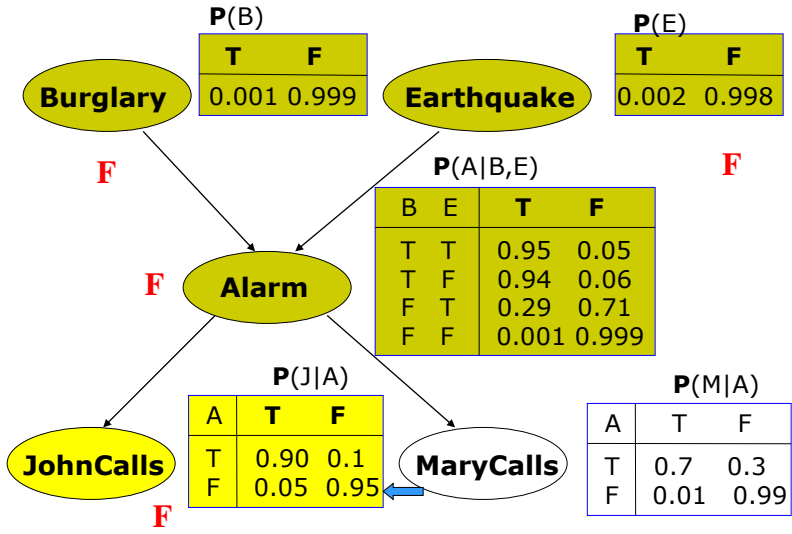
CS 3750 Advanced Machine Learning

## BBN sampling example



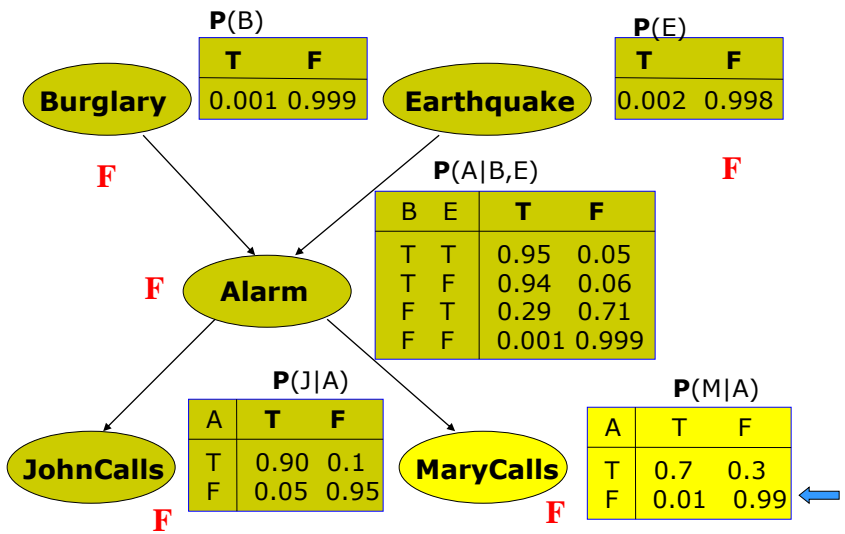
CS 3750 Advanced Machine Learning

## BBN sampling example



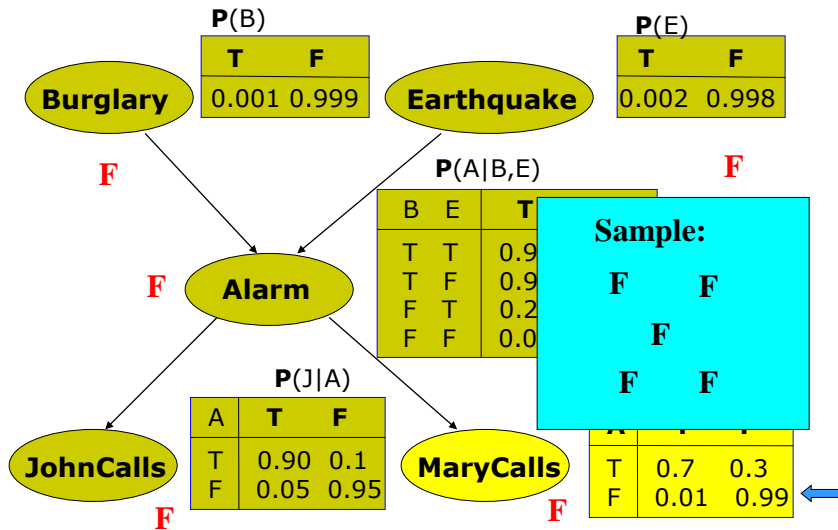
CS 3750 Advanced Machine Learning

## BBN sampling example



CS 3750 Advanced Machine Learning

## BBN sampling example



CS 3750 Advanced Machine Learning

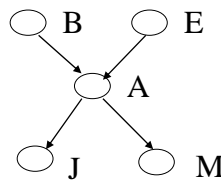
## Monte Carlo inference: BBNs

**Challenge 1: How to generate  $M$  (unbiased) examples from the target distribution  $P(X)$  defined by BBN?**

- **Good news: Sample generation for the full joint defined by the BBN is easy**

- One top down sweep through the network lets us generate one example according to  $P(X)$

- **Example:**



Examples are generated in a top down manner, following the links

- Repeat many times to get enough of examples

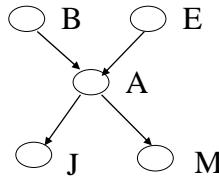
CS 3750 Advanced Machine Learning



## Monte Carlo inference: BBNs

Knowing how to generate efficiently examples from the full joint lets us efficiently estimate:

- Joint probabilities over a subset variables
- Marginals on variables
- **Example:**



The probability is approximated using sample frequency

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

← # samples with  $B = T, J = T$   
 ← total # samples

CS 3750 Advanced Machine Learning

## Monte Carlo inference: BBNs

- **MC approximation of conditional probabilities:**
  - The probability can approximated using sample frequencies
  - **Example:**

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

← # samples with  $B = T, J = T$   
 ← # samples with  $J = T$

- **Solution 1 (rejection sampling):**
  - Generate examples from  $P(X)$  which we know how to do efficiently
  - Use only samples that agree with the condition ( $J=T$ ), the remaining samples are rejected
- **Problem:** many examples are rejected. What if  $P(J=T)$  is very small?

CS 3750 Advanced Machine Learning

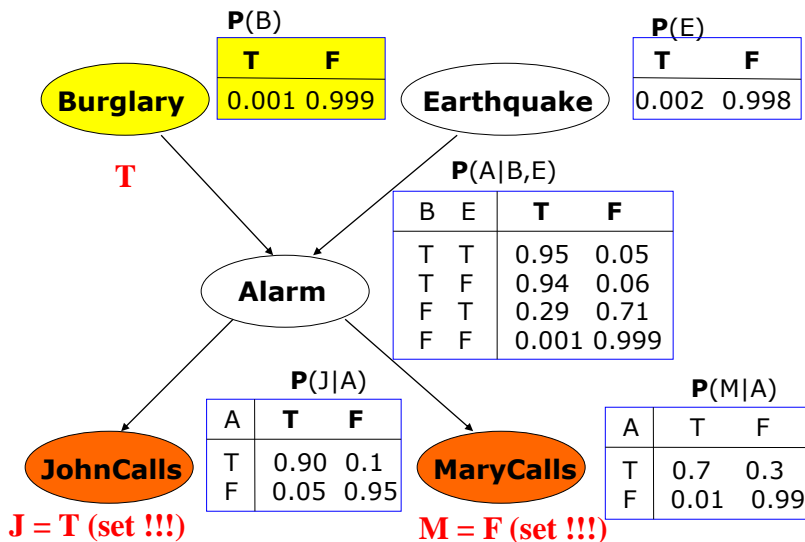
## Monte Carlo inference: BBNs

- MC approximation of conditional probabilities
- Solution 2 (likelihood weighting)
  - Avoids inefficiencies of rejection sampling
  - **Idea:** generate only samples consistent with an evidence (or conditioning event); If the value is set no sampling
- **Problem:** using simple counts is not enough since these may occur with different probabilities
- Likelihood weighting:
  - **With every sample keep a weight with which it should count towards the estimate**

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

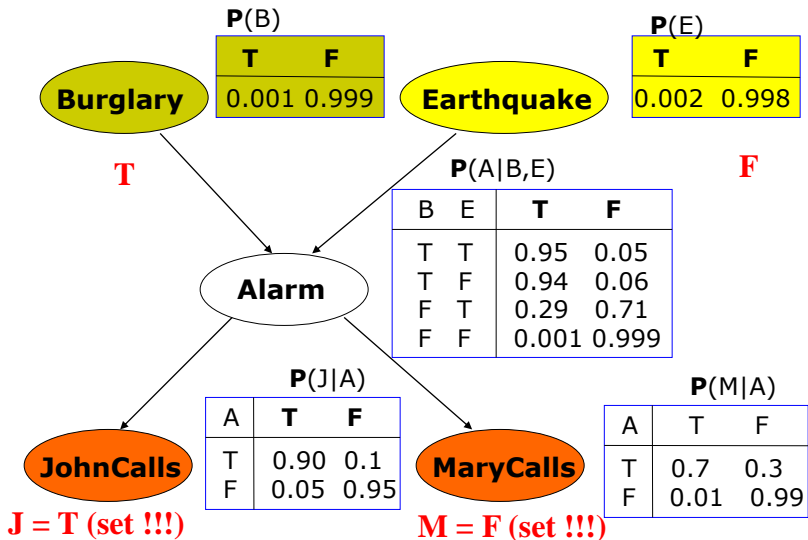
CS 3750 Advanced Machine Learning

## BBN likelihood weighting example



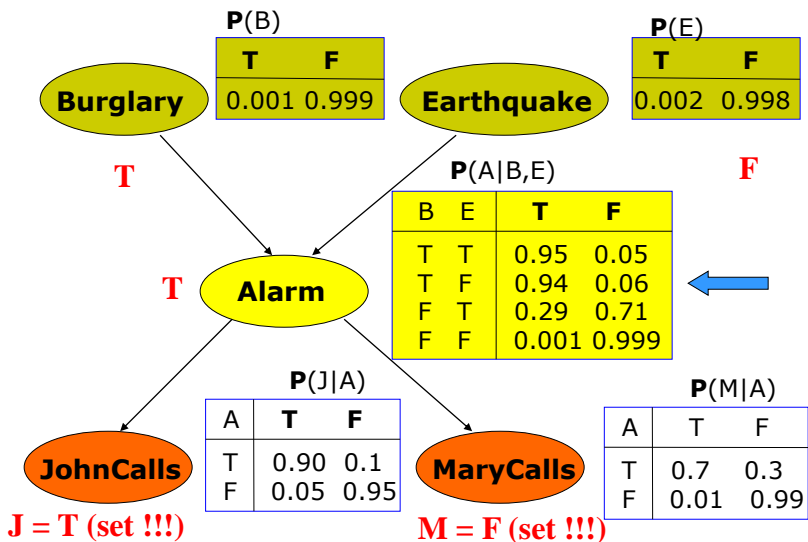
CS 3750 Advanced Machine Learning

## BBN likelihood weighting example



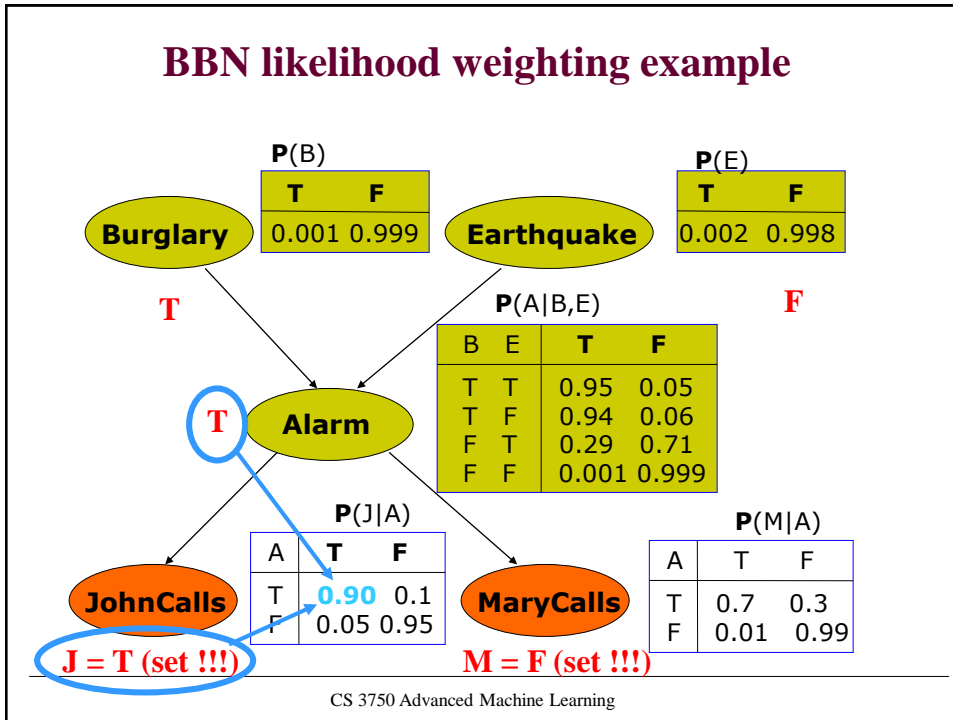
CS 3750 Advanced Machine Learning

## BBN likelihood weighting example

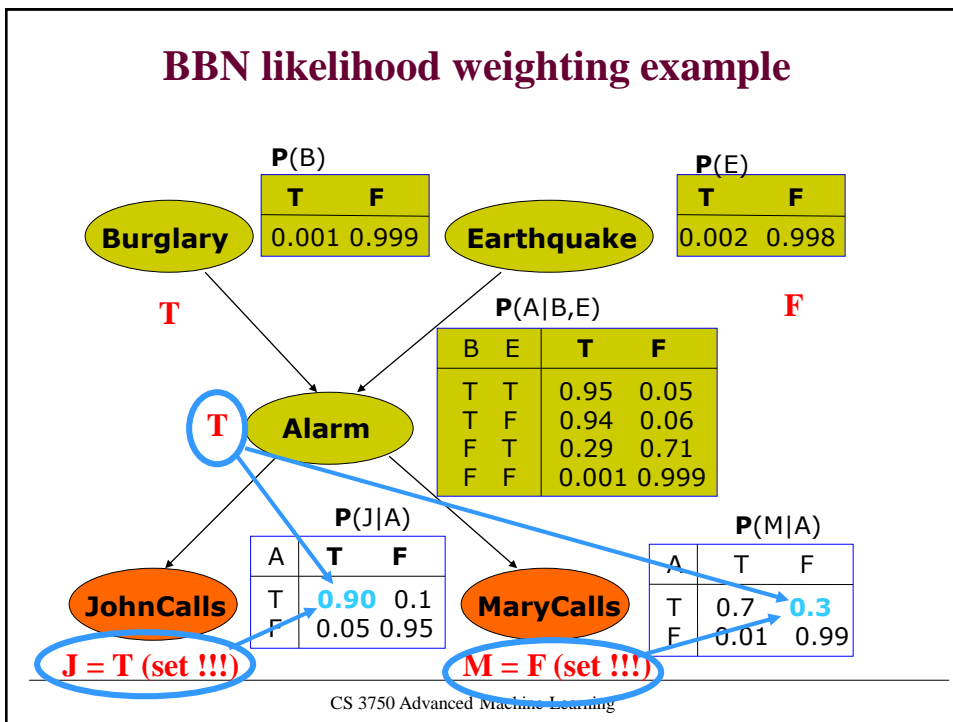


CS 3750 Advanced Machine Learning

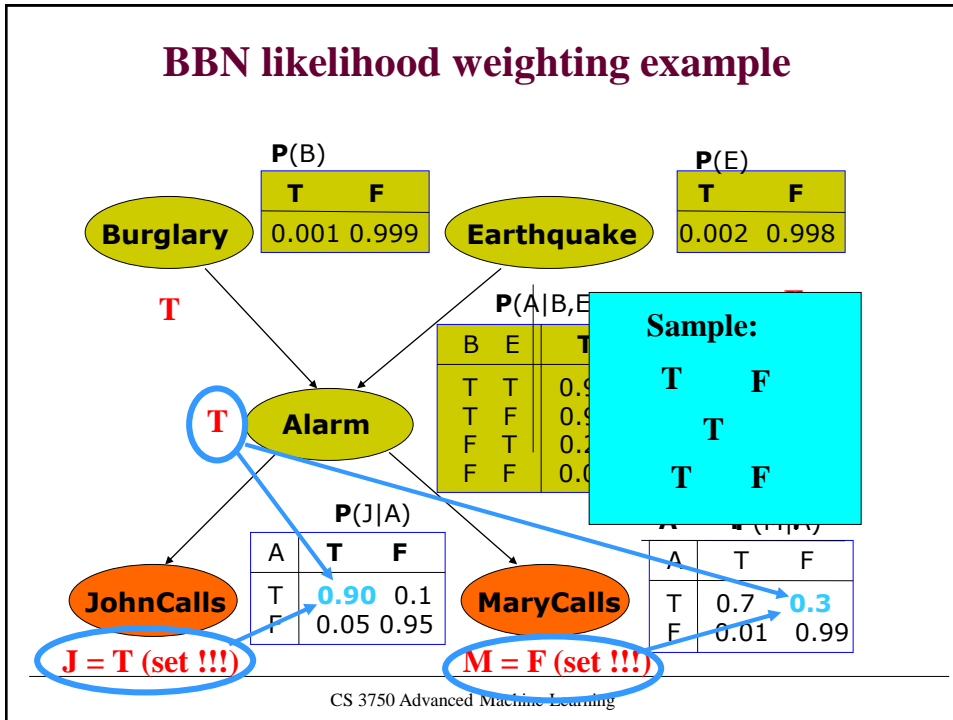
## BBN likelihood weighting example



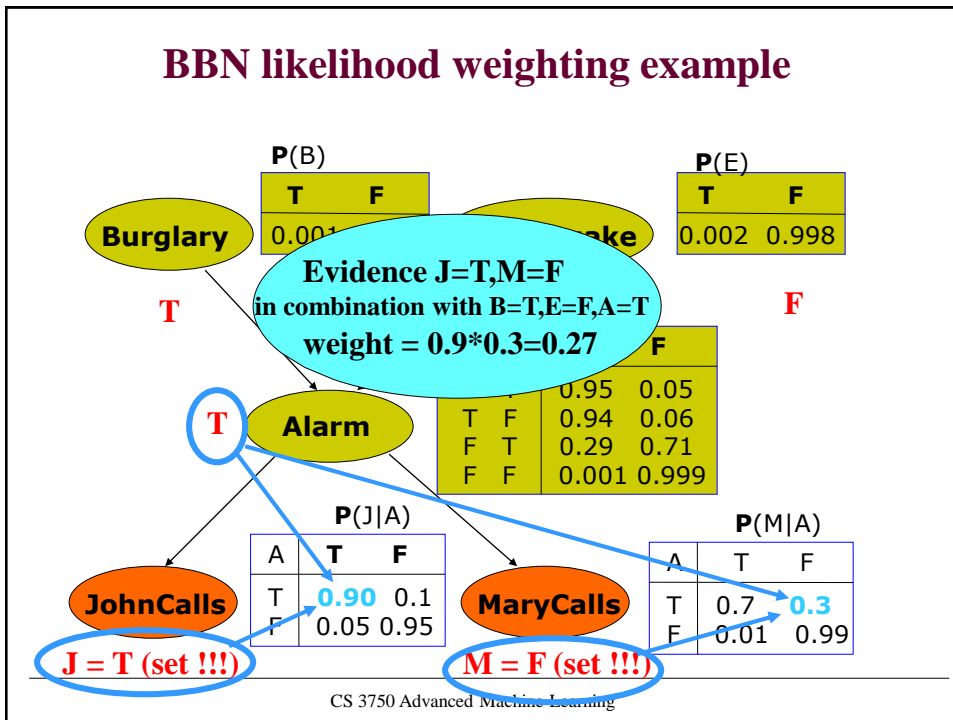
## BBN likelihood weighting example



## BBN likelihood weighting example

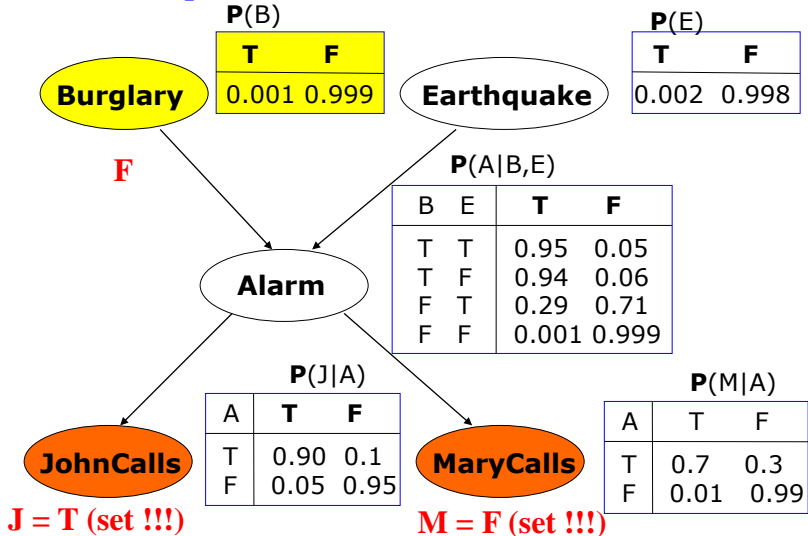


## BBN likelihood weighting example



## BBN likelihood weighting example

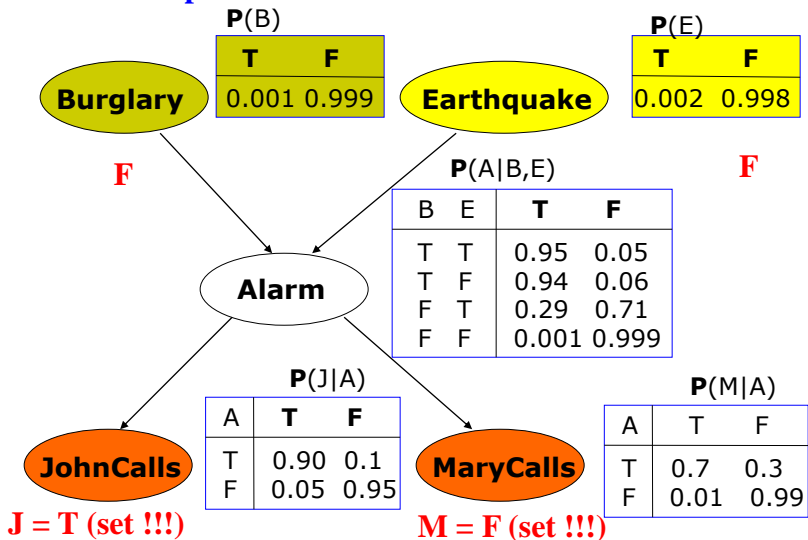
Second sample



CS 3750 Advanced Machine Learning

## BBN likelihood weighting example

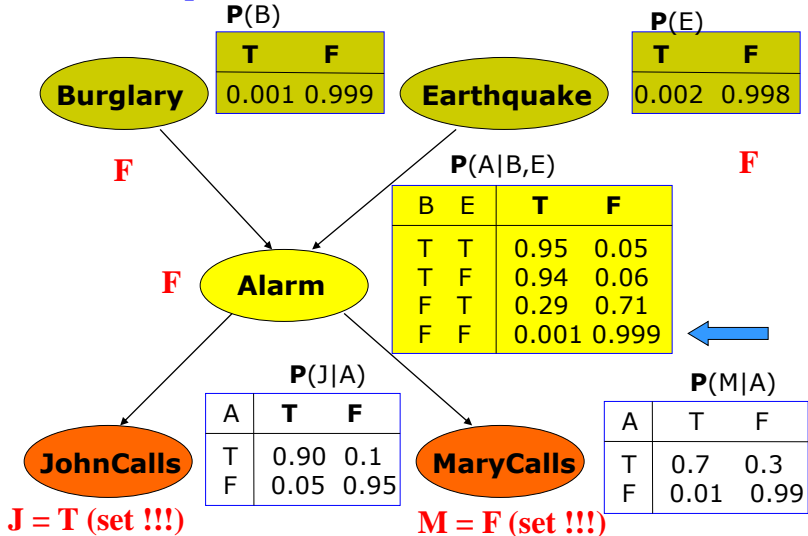
Second sample



CS 3750 Advanced Machine Learning

## BBN likelihood weighting example

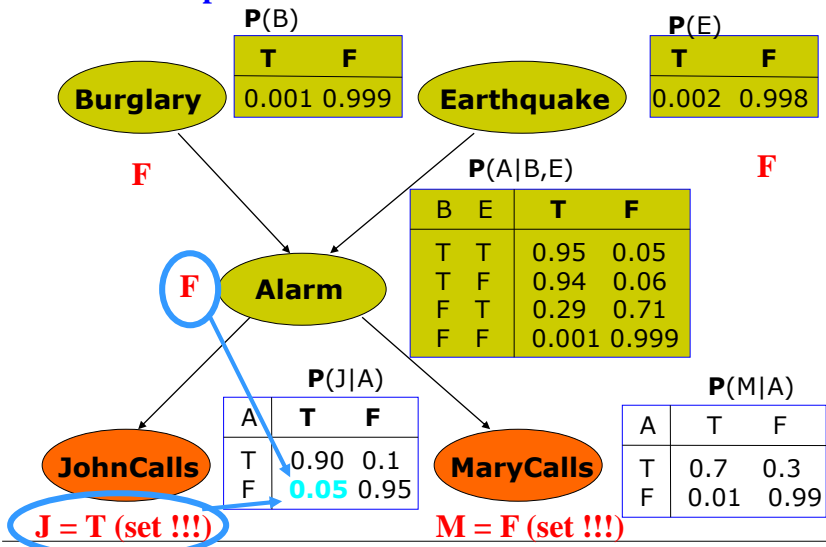
Second sample



CS 3750 Advanced Machine Learning

## BBN likelihood weighting example

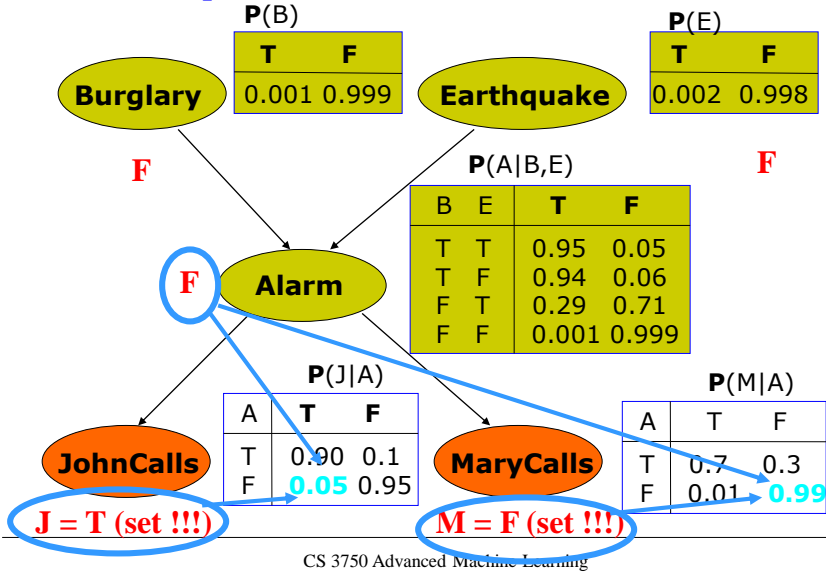
Second sample



CS 3750 Advanced Machine Learning

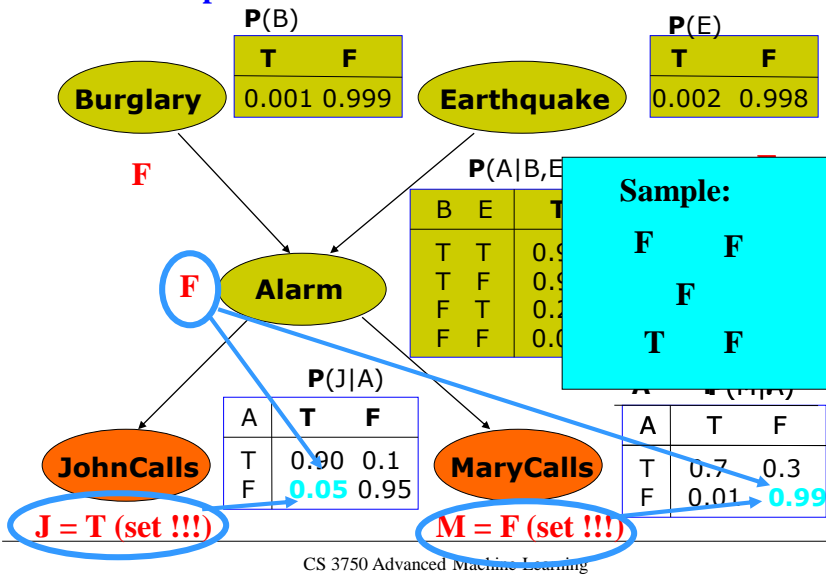
## BBN likelihood weighting example

Second sample



## BBN likelihood weighting example

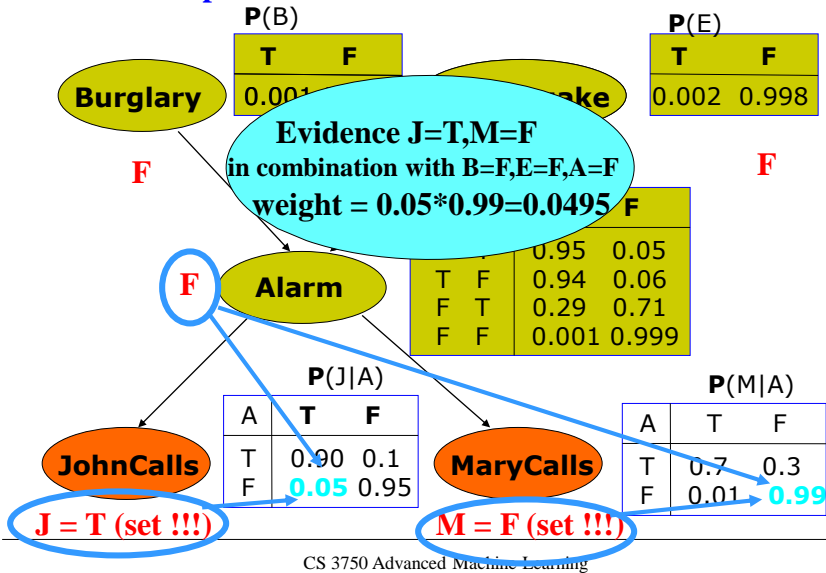
Second sample





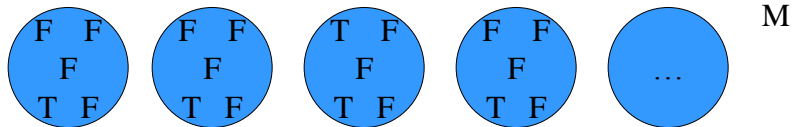
## BBN likelihood weighting example

Second sample



## Likelihood weighting

- Assume we have generated the following  $M$  samples:



- If we calculate the estimate:

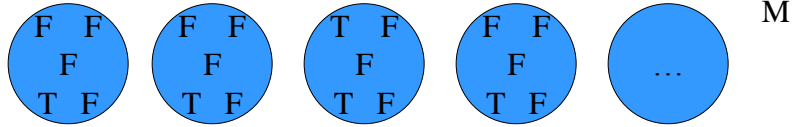
$$P(B=T | J=T, M=F) = \frac{\# \text{sample\_with}(B=T)}{\# \text{total\_sample}}$$

a less likely sample from  $P(X)$  may be generated more often.

- For example, sample is generated more often than in  $P(X)$
- So the samples are not consistent with  $P(X)$ .

## 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  $P(e)$ .



35

## Likelihood weighting

- How to compute weights for the sample?
- Assume the query  $P(B = T \mid J = T, M = F)$
- Likelihood weighting:
  - With every sample keep a weight with which it should count towards the estimate

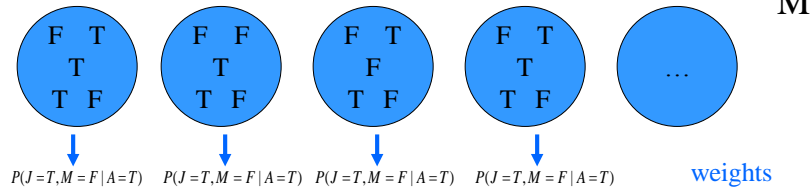
$$\tilde{P}(B = T \mid J = T, M = F) = \frac{\sum_{i=1}^M \mathbb{1}\{B^{(i)} = T\} w^{(i)}}{\sum_{i=1}^M w^{(i)}}$$

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

36

## Likelihood weighting

- Assume  $M$  samples where evidence is enforced:



- We can use  $P(e)$  to weight each sample and correct the bias.
- The correct estimate is then:

$$\tilde{P}(A=T | J=T, M=F) = \frac{\sum_{i=1}^M \mathbb{I}\{A^{(i)}=T\} w^{(i)}}{\sum_{i=1}^M w^{(i)}}$$

37

## Monte Carlo inference: MRFs

**Challenge:** How to generate  $M$  (unbiased) examples from the target distribution  $P(X)$  defined by an MRF?

- Trivial solution:**
  - calculate and store the probability of each configuration
  - Pick randomly a configuration based on its probability
- Problem:** terribly inefficient for a large number of variables
- Can we do better, similarly to BBN?
- In general, sampling  $P(X)$  or  $P(X'|Evidence)$  can be hard?

**Next:** avoid sampling  $P(X)$  by sampling  $Q(X)$

38

## Importance Sampling

- An approach for estimating the expectation of a function  $f(x)$  relative to some distribution  $P(X)$  (**target distribution**)
- generally, we can estimate this expectation by generating samples  $x[1], \dots, x[M]$  from  $P$ , and then estimating

$$E_p[f] = \frac{1}{M} \sum_{m=1}^M f(x[m])$$

- However, we might prefer to generate samples from a different distribution  $Q$  (**proposal or sampling distribution**) instead, since it might be impossible or computationally very expensive to generate samples directly from  $P(X)$ .
- $Q$  can be arbitrary, but it should dominate  $P$ , i.e.  $Q(x) > 0$  whenever  $P(x) > 0$

CS 3750 Advanced Machine Learning

## Unnormalized Importance Sampling

- Since we generate samples from  $Q$  instead of  $P$ ,
- we need to adjust our estimator to compensate for the incorrect sampling distribution.

$$E_{p(x)}[f(X)] = E_{Q(x)}\left[f(x) \frac{P(x)}{Q(x)}\right]$$

- So we can use standard estimator for expectations relative to  $Q$ .
- **Method:** We generate a set of  $M$  samples  $D = \{x[1], \dots, x[M]\}$  from  $Q$ , and estimate:

$$\hat{E}_D(f) = \frac{1}{M} \sum_{m=1}^M f(x[m]) \frac{P(x[m])}{Q(x[m])}$$

CS 3750 Advanced Machine Learning

## Importance sampling

- This is an unbiased estimator: its mean for any data set is precisely the desired value

$$w(x) = P(x)/Q(x) \quad \text{- a weighting function, or a correction weight}$$

- We can estimate the distribution of the estimator around its mean: as  $M \rightarrow \infty$

$$E_{Q(X)}[f(X)w(X)] - E_{P(X)}[f(X)] \propto N(0; \sigma_Q^2 / M)$$

$$\text{where } \sigma_Q^2 = [E_{Q(X)}[(f(X)w(X))^2]] - (E_{Q(X)}[f(X)w(X)])^2$$

$$\sigma_Q^2 = [E_{Q(X)}[(f(X)w(X))^2]] - (E_{P(X)}[f(X)])^2$$

---

CS 3750 Advanced Machine Learning

## Importance sampling

- When  $f(X)=1$ , the variance is simply the variance of the weighting function  $P(X)/Q(X)$ . Thus, the more different  $Q$  is from  $P$ , the higher is the variance of the estimator.
- In general, the lowest variance is achieved when

$$Q(X) \propto |f(X)| P(X)$$

- We should avoid cases where our sampling probability  $Q(X) \ll P(X)f(X)$  in any part of the space, as these cases can lead to very large or even infinite variance.
- Problem with un-normalized IS:  $P$  is assumed to be known

---

CS 3750 Advanced Machine Learning

## Normalized Importance Sampling

- When  $P$  is only known up to a normalizing constant  $\alpha$
- We have access to a function  $P'(X)$ , such that  $P'$  is not a normalized distribution, but  $P'(X) = \alpha P(X)$
- In this context, we cannot define the weights relative to  $P$ , so we define:

$$w(X) = \frac{P'(X)}{Q(X)}$$

$$\begin{aligned} E_{P(X)}[f(X)] &= \sum_x P(x) f(x) = \sum_x Q(x) f(x) \frac{P(X)}{Q(x)} = \frac{1}{\alpha} \sum_x Q(x) f(x) \frac{P'(x)}{Q(x)} \\ &= \frac{1}{\alpha} E_{Q(x)}[f(X)w(X)] = \frac{E_{Q(x)}[f(X)w(X)]}{E_{Q(x)}[w(X)]} \end{aligned}$$

Why? 
$$E_{Q(x)}[w(X)] = \sum_x Q(x) \frac{P'(x)}{Q(x)} = \sum_x P'(x) = \alpha$$

CS 3750 Advanced Machine Learning

## Importance sampling

- Using an empirical estimator for both the numerator and denominator, we can estimate:

$$\hat{E}_D(f) = \frac{\sum_{m=1}^M f(x[m])w(x[m])}{\sum_{m=1}^M w(x[m])}$$

- Although the normalized estimator is biased, its variance is typically lower than that of the unnormalized estimator. This reduction in variance often outweighs the bias term.
- So normalized estimator is often used in place of the unnormalized estimator, even in cases where  $P$  is known and we can sample from it effectively.

CS 3750 Advanced Machine Learning

## Importance sampling for estimating conditional probabilities in BBNs

Assume a Bayesian Network

- We want to calculate  $P(x'|evidence)$
- This is hard if we need to go opposite the links and account for the effect of evidence on non-descendants

**Objective:** generate samples efficiently using a simpler proposal distribution  $Q(x)$

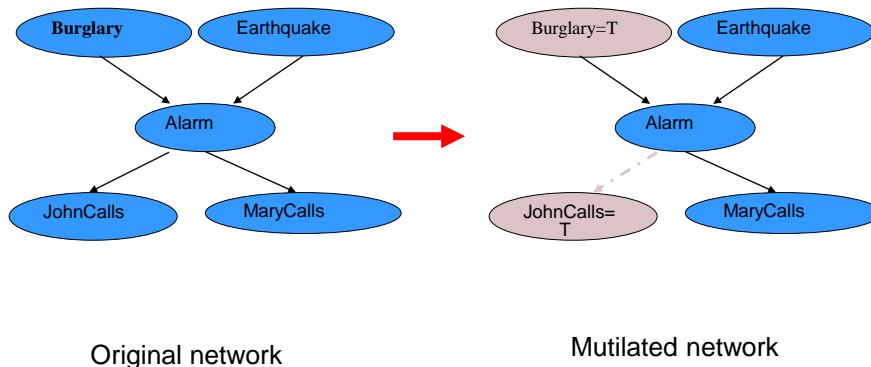
**Solution:** a **mutilated belief network** (Koller, Friedman 2009)

- **Idea:**
  - Avoid propagation of evidence effects to nondescendants;
  - Disconnect all variables in the evidence from their parents

45

## Mutilated Belief network

- Assume we want to calculate  $P(x | B=T, J=T)$  in the Alarm network
- Use  $B=T$  and  $J=T$  to build a mutilated network



46

## Mutilated Belief network

- Assume the evidence is  $J=j^*$  and  $B=b^*$

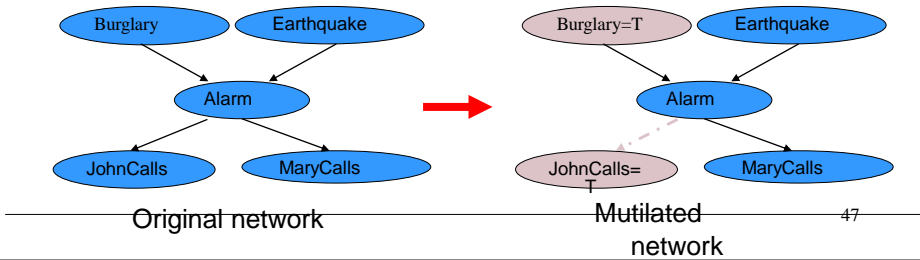
- Original network:

$$P(E=e, A=a, M=m, J=j^*, B=b^*) = P(b^*)P(e)P(a|b^*, e)P(j^*|a)P(m|a)$$

- Mutilated network:

$$Q(E=e, A=a, M=m, J=j^*, B=b^*) = P(e)P(a|b^*, e)P(m|a)$$

- Note that  $w(x) = \frac{P(x)}{Q(x)} = P(b^*)P(j^*|a)$



## Mutilated Belief network

- Assume the evidence is  $J=j^*$  and  $B=b^*$

- Original network:

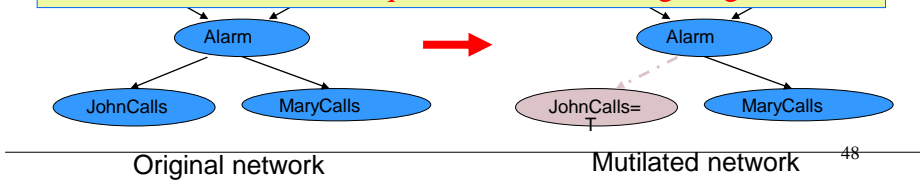
$$P(E=e, A=a, M=m, J=j^*, B=b^*) = P(b^*)P(e)P(a|b^*, e)P(j^*|a)P(m|a)$$

- Mutilated network:

$$Q(E=e, A=a, M=m, J=j^*, B=b^*) = P(e)P(a|b^*, e)P(m|a)$$

- Note that  $w(x) = \frac{P(x)}{Q(x)} = P(b^*)P(j^*|a)$

So importance sampling with a proposal distribution based on mutilated network is equal to likelihood weighting





## Likelihood Weighting

- **Question:** When to stop? How many samples do we need to see?
- **Intuition:** not every sample contribute equally to the quality of the estimate. A sample with a high weight is more compatible with the evidence  $e$ , and may provide us with more information.
- **Solution:** We stop sampling when the total weight of the generated samples reaches a pre-defined value.
- **Benefits:** It allows early stopping in cases where we were lucky in our random choice of samples.

49

## Markov chain Monte Carlo

- **Likelihood weighting:** samples are generated according to  $Q$  and every sample from  $Q$  is reweighted according to its likelihood, but the  $Q$  distribution may be very far from the target
- **MCMC** is a strategy for generating samples from the target distribution, including conditional distributions
- **MCMC:**
  - Markov chain defines a sampling process that
  - initially generates samples very different from the target distribution (e.g. posterior)
  - but gradually refines the samples so that they are closer and closer to the posterior.

CS 3750 Advanced Machine Learning

## MCMC

- The construction of a Markov chain requires two basic ingredients
  - a transition matrix  $P$
  - an initial distribution  $\pi_0$
- Assume a finite set  $S = \{1, \dots, m\}$  of states, then a **transition matrix** is

$$P = \begin{pmatrix} p_{11} & p_{12} & \cdots & p_{1m} \\ p_{21} & p_{22} & \cdots & p_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ p_{m1} & p_{m2} & \cdots & p_{mm} \end{pmatrix}$$

Where  $p_{ij} \geq 0 \quad \forall (i, j) \in S^2$  and  $\sum_{j \in S} p_{ij} = 1 \quad \forall i \in S$

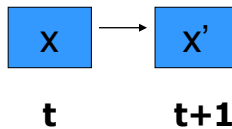
## Markov Chain

- Markov chain defines a random process of selecting states

$x^{(0)}, x^{(1)}, \dots, x^{(m)}, \dots$

Initial state selected based on  $\pi_0$

Subsequent states selected based on the previous state and the transition matrix



- Chain Dynamics**

$$P^{(t+1)}(X^{(t+1)} = x') = \sum_{x \in \text{Dom}(X)} P^{(t)}(X^{(t)} = x) T(x \rightarrow x')$$

Probability of a state  $x'$  being selected at time  $t+1$

transition matrix

## MCMC

- **Markov chain** satisfies

$$P(X_{n+1} = j | X_0 = i_0, X_1 = i_1, \dots, X_n = i_n) = P(X_{n+1} = j | X_n = i_n)$$

- **Irreducibility:** A MC is called irreducible (or un-decomposable) if there is a positive transition probability for all pairs of states within a limited number of steps
- In irreducible chains there may still exist a periodic structure such that for each state  $i \in \mathcal{S}$ , the set of possible return times to  $i$  when starting in  $i$  is a subset of the set  $p\mathbb{N} = \{p, 2p, 3p, \dots\}$  containing all but a finite set of these elements. The smallest number  $p$  with this property is the so-called **period of the chain**

$$p = \gcd\{n \in \mathbb{N} : p_{ii}^{(n)} > 0\}$$

---

CS 3750 Advanced Machine Learning

## MCMC

- **Aperiodicity:** An irreducible chain is called aperiodic (or acyclic) if the period  $p$  equals 1 or, equivalently, if for all pairs of states there is an integer  $n_{ij}$  such that for all  $n \geq n_{ij}$ , the probability  $p_{ij}^{(n)} > 0$ .
- If a Markov chain satisfy both **irreducibility and aperiodicity**, then it **converges to an invariant distribution  $q(x)$**
- A Markov chain with transition matrix  $P$  will have an equilibrium distribution  $q$  **iff**  $q = qP$ .
- A sufficient, but not necessary, condition to ensure a particular  $q(x)$  is the invariant distribution of transition matrix  $P$  is the following **reversibility (detailed balance) condition**

$$q(x^i)P(x^{i-1} | x^i) = q(x^{i-1})P(x^i | x^{i-1})$$

---

CS 3750 Advanced Machine Learning

## Markov Chain Monte Carlo

**Objective:** generate samples from the target distribution (e.g. posterior)

- **Idea:**

Markov chain defines a sampling process that:

- initially generates samples very different from the target posterior
- but gradually refines the samples so that they are closer and closer to the target distribution