

# Probabilistic Graphical Models

## Lecture 12: Approximate Inference: Sampling

Matthew Walter

TTI-Chicago

May 19, 2020

Some slide content courtesy of Eric Xing

# Approaches to Inference

- Exact inference algorithms
  - Variable elimination
  - Belief propagation
    - Sum-product message passing
    - Belief update message passing
- Approximate inference algorithms
  - Variational methods
    - Loopy belief propagation
    - Mean-field approximation
  - Sample-based inference: Markov chain Monte Carlo (MCMC)

# Monte Carlo Methods

- Suppose that we have a distribution  $P$  over  $\mathcal{X}$
- We want to estimate  $\mathbb{E}[f(\mathcal{X})]$  for some function  $f(\mathcal{X})$
- Generate  $M$  samples (“particles”) from  $P$   $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(M)}\}$  and estimate the expectation of  $f(\mathcal{X})$

$$\hat{\mathbb{E}}[f(\mathbf{X})] = \frac{1}{M} \sum_{m=1}^M f(\mathbf{x}^{(m)})$$

- Suppose that we want the marginal  $P(\mathbf{Y} = \mathbf{y})$  for some  $\mathbf{Y} \subseteq \mathcal{X}$

# Monte Carlo Methods

- Suppose that we have a distribution  $P$  over  $\mathcal{X}$
- We want to estimate  $\mathbb{E}[f(\mathcal{X})]$  for some function  $f(\mathcal{X})$
- Generate  $M$  samples (“particles”) from  $P$   $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(M)}\}$  and estimate the expectation of  $f(\mathcal{X})$

$$\hat{\mathbb{E}}[f(\mathbf{X})] = \frac{1}{M} \sum_{m=1}^M f(\mathbf{x}^{(m)})$$

- Suppose that we want the marginal  $P(\mathbf{Y} = \mathbf{y})$  for some  $\mathbf{Y} \subseteq \mathcal{X}$   
 $\Rightarrow$  If  $f(\mathcal{X}) = \mathbb{1}(\xi\langle \mathbf{Y} \rangle = \mathbf{y})$ , this gives us the marginal  $P(\mathbf{y})$

$$\hat{P}(\mathbf{y}) = \frac{1}{M} \sum_{m=1}^M \mathbb{1}(\langle \mathbf{Y} \rangle = \mathbf{y})$$

(i.e., the fraction of particles consistent with  $\mathbf{Y} = \mathbf{y}$ )

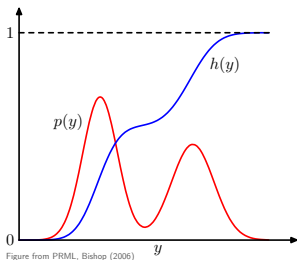
- Draw random samples from the desired distribution
- Yields a stochastic representation of a complex distribution
- Exhibits *asymptotic* convergence
- Questions:
  - How do we sample from a given distribution (i.e., not all distributions can be easily sampled)?
  - How can we make better use of samples (not all samples are equal)?
  - How do we know when we've sampled enough?

# Sampling from Distributions

- Key idea: Convert this to sampling from a uniform distribution (easy)
- Consider a discrete distribution, i.e., a multinomial  $P(X)$  for  $\text{Val}(X) = \{x^1, \dots, x^k\}$  with  $P(X = x^i) = \theta_i$ 
  - 1 Partition interval  $[0, 1]$  into  $k$  subintervals  $[0, \theta_1), [\theta_1, \theta_1 + \theta_2), \dots]$
  - 2 Sample  $s \sim \text{Uniform}[0, 1]$  and choose  $x^i$  if  $s$  is in the  $i^{\text{th}}$  interval

# Sampling from Distributions

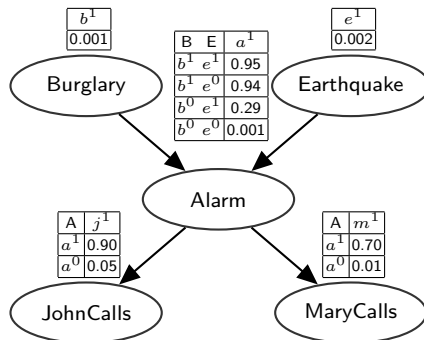
- Key idea: Convert this to sampling from a uniform distribution (easy)
- Consider a discrete distribution, i.e., a multinomial  $P(X)$  for  $\text{Val}(X) = \{x^1, \dots, x^k\}$  with  $P(X = x^i) = \theta_i$ 
  - 1 Partition interval  $[0, 1]$  into  $k$  subintervals  $[0, \theta_1), [\theta_1, \theta_1 + \theta_2), \dots]$
  - 2 Sample  $s \sim \text{Uniform}[0, 1]$  and choose  $x^i$  if  $s$  is in the  $i^{\text{th}}$  interval
- Similarly, for continuous distributions  $p(x)$  we consider the cumulative distribution  $h(y) = \int_{-\infty}^y p(x) dx$ : Sample  $s \sim \text{Uniform}[0, 1]$  and return  $h^{-1}(s)$  (not always easy: we can't always compute and invert  $h(y)$ )



# Forward Sampling

- Arguably the simplest approach to generating samples (particles)
- Involves *directly* sampling from  $P(\mathcal{X})$
- Requires knowledge of the partition function for MRFs
- Primarily restricted to Bayesian networks
- Procedure: Given a ordering of the random variables,
  - 1 Sample each  $X_i$  according to its CPD  $P(X_i | \text{Pa}_{X_i}^G)$  using the current sampled values for its parents
  - 2 Repeat  $M$  times







## Forward Sampling: Example

**Procedure:** Given a ordering of the random variables,

- 1 Sample each  $X_i$  according to it's CPD  $P(X_i | \text{Pa}_{X_i}^G)$  using the current sampled values for its parents
- 2 Repeat  $M$  times

**Problem:** It is difficult to get sufficient samples of rare events

- What about  $P(J | A^1)$ ? We only have one sample:

$$P(J | a^1) = \frac{P(J, a^1)}{P(a^1)} = \{0, 1\}$$

- What about  $P(J | b^1)$ ? We don't have any samples!

E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E0	B0	A0	M0	J1
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E1	B0	A1	M1	J1
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0

# How Many Samples Are Necessary?

- Suppose that we are interested in a particular event  $\mathbf{Y} = \mathbf{y}$
- $\mathbb{1}(\langle \mathbf{Y} \rangle = \mathbf{y})$  is a Bernoulli random variable with success  $P(\mathbf{y})$ 
  - $M$  samples define  $M$  independent Bernoulli trials

# How Many Samples Are Necessary?

- Suppose that we are interested in a particular event  $\mathbf{Y} = \mathbf{y}$
- $\mathbb{1}(\langle \mathbf{Y} \rangle = \mathbf{y})$  is a Bernoulli random variable with success  $P(\mathbf{y})$ 
  - $M$  samples define  $M$  independent Bernoulli trials
- **Hoeffding bound:** Consider a sequence  $\mathcal{D} = \{x^{(1)}, x^{(2)}, \dots, x^{(M)}\}$  of  $M$  independent Bernoulli trials with probability  $p$ . Letting  $T_{\mathcal{D}} = \frac{1}{M} \sum_m x^{(m)}$  (e.g., fraction of heads), for some  $\epsilon > 0$

$$P_{\mathcal{D}}(T_{\mathcal{D}} > p + \epsilon) \leq e^{-2M\epsilon^2}$$

$$P_{\mathcal{D}}(T_{\mathcal{D}} < p - \epsilon) \leq e^{-2M\epsilon^2}$$

# How Many Samples Are Necessary?

- Suppose that we are interested in a particular event  $\mathbf{Y} = \mathbf{y}$
- $\mathbb{1}(\langle \mathbf{Y} \rangle = \mathbf{y})$  is a Bernoulli random variable with success  $P(\mathbf{y})$ 
  - $M$  samples define  $M$  independent Bernoulli trials
- **Hoeffding bound:** Consider a sequence  $\mathcal{D} = \{x^{(1)}, x^{(2)}, \dots, x^{(M)}\}$  of  $M$  independent Bernoulli trials with probability  $p$ . Letting  $T_{\mathcal{D}} = \frac{1}{M} \sum_m x^{(m)}$  (e.g., fraction of heads), for some  $\epsilon > 0$

$$P_{\mathcal{D}}(T_{\mathcal{D}} > p + \epsilon) \leq e^{-2M\epsilon^2}$$

$$P_{\mathcal{D}}(T_{\mathcal{D}} < p - \epsilon) \leq e^{-2M\epsilon^2}$$

- Using the Hoeffding bound, we get

$$P_{\mathcal{D}}(\hat{P}_{\mathcal{D}} \notin [P(\mathbf{y}) - \epsilon, P(\mathbf{y}) + \epsilon]) \leq 2e^{-2M\epsilon^2}$$

# How Many Samples Are Necessary?

- Suppose that we are interested in a particular event  $\mathbf{Y} = \mathbf{y}$
- $\mathbb{1}(\langle \mathbf{Y} \rangle = \mathbf{y})$  is a Bernoulli random variable with success  $P(\mathbf{y})$ 
  - $M$  samples define  $M$  independent Bernoulli trials
- **Hoeffding bound:** Consider a sequence  $\mathcal{D} = \{x^{(1)}, x^{(2)}, \dots, x^{(M)}\}$  of  $M$  independent Bernoulli trials with probability  $p$ . Letting  $T_{\mathcal{D}} = \frac{1}{M} \sum_m x^{(m)}$  (e.g., fraction of heads), for some  $\epsilon > 0$

$$P_{\mathcal{D}}(T_{\mathcal{D}} > p + \epsilon) \leq e^{-2M\epsilon^2}$$

$$P_{\mathcal{D}}(T_{\mathcal{D}} < p - \epsilon) \leq e^{-2M\epsilon^2}$$

- Using the Hoeffding bound, we get

$$P_{\mathcal{D}}(\hat{P}_{\mathcal{D}} \notin [P(\mathbf{y}) - \epsilon, P(\mathbf{y}) + \epsilon]) \leq 2e^{-2M\epsilon^2}$$

- Thus, to achieve an estimate within  $\epsilon$  of the true marginal (absolute error) with probability at least  $1 - \delta$

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

# How Many Samples Are Necessary?

- **Chernoff bound:** Consider a sequence  $\mathcal{D} = \{x^{(1)}, x^{(2)}, \dots, x^{(M)}\}$  of  $M$  independent Bernoulli trials with probability  $p$ . For  $T_{\mathcal{D}} = \frac{1}{M} \sum_m x^{(m)}$  we have

$$P_{\mathcal{D}}(T_{\mathcal{D}} > p(1 + \epsilon)) \leq e^{-Mp\epsilon^2/3}$$

$$P_{\mathcal{D}}(T_{\mathcal{D}} < p(1 - \epsilon)) \leq e^{-Mp\epsilon^2/3}$$

- Using the Chernoff bound, we can compute the likelihood associated with the *relative error* (suitable for disparate likelihoods)

$$P_{\mathcal{D}}\left(\hat{P}_{\mathcal{D}} \notin P(\mathbf{y})(1 \pm \epsilon)\right) \leq 2e^{-MP(\mathbf{y})\epsilon^2/3}$$



# How Many Samples Are Necessary?

- **Chernoff bound:** Consider a sequence  $\mathcal{D} = \{x^{(1)}, x^{(2)}, \dots, x^{(M)}\}$  of  $M$  independent Bernoulli trials with probability  $p$ . For  $T_{\mathcal{D}} = \frac{1}{M} \sum_m x^{(m)}$  we have

$$P_{\mathcal{D}}(T_{\mathcal{D}} > p(1 + \epsilon)) \leq e^{-Mp\epsilon^2/3}$$

$$P_{\mathcal{D}}(T_{\mathcal{D}} < p(1 - \epsilon)) \leq e^{-Mp\epsilon^2/3}$$

- Using the Chernoff bound, we can compute the likelihood associated with the *relative error* (suitable for disparate likelihoods)

$$P_{\mathcal{D}}\left(\hat{P}_{\mathcal{D}} \notin P(\mathbf{y})(1 \pm \epsilon)\right) \leq 2e^{-MP(\mathbf{y})\epsilon^2/3}$$

- Thus, to achieve a relative error no more than  $\epsilon$  with probability  $1 - \delta$

$$M \geq 3 \frac{\ln(2/\delta)}{P(\mathbf{y})\epsilon^2}$$

# How Many Samples Are Necessary?

- **Chernoff bound:** Consider a sequence  $\mathcal{D} = \{x^{(1)}, x^{(2)}, \dots, x^{(M)}\}$  of  $M$  independent Bernoulli trials with probability  $p$ . For  $T_{\mathcal{D}} = \frac{1}{M} \sum_m x^{(m)}$  we have

$$P_{\mathcal{D}}(T_{\mathcal{D}} > p(1 + \epsilon)) \leq e^{-Mp\epsilon^2/3}$$

$$P_{\mathcal{D}}(T_{\mathcal{D}} < p(1 - \epsilon)) \leq e^{-Mp\epsilon^2/3}$$

- Using the Chernoff bound, we can compute the likelihood associated with the *relative error* (suitable for disparate likelihoods)

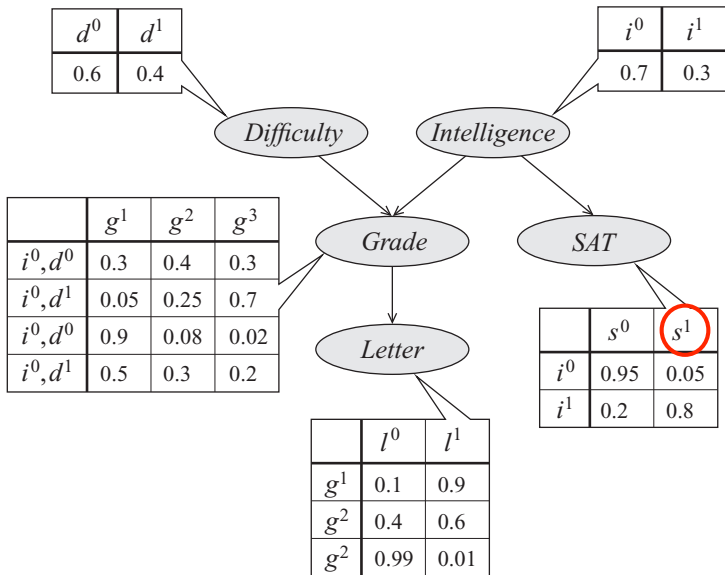
$$P_{\mathcal{D}}\left(\hat{P}_{\mathcal{D}} \notin P(\mathbf{y})(1 \pm \epsilon)\right) \leq 2e^{-MP(\mathbf{y})\epsilon^2/3}$$

- Thus, to achieve a relative error no more than  $\epsilon$  with probability  $1 - \delta$

$$M \geq 3 \frac{\ln(2/\delta)}{P(\mathbf{y})\epsilon^2}$$

- The number of samples is inversely proportional to  $P(\mathbf{y})$ 
  - Low probability events require more samples
  - Oftentimes, we don't know  $P(\mathbf{y})$

# Dealing with Evidence: Example



# Dealing with Evidence

- How can we use sampling to estimate  $P(\mathbf{y} \mid \mathbf{E} = \mathbf{e})$
- One approach is to sample from  $P(\mathbf{y} \mid \mathbf{E} = \mathbf{e})$  via *rejection sampling*
  - Draw  $\tilde{M}$  samples  $\mathbf{x} \sim P(\mathbf{X})$  as before
  - Discard (reject) samples that are inconsistent with  $\mathbf{E} = \mathbf{e}$
  - Resulting samples are from the posterior  $P(\mathbf{X} \mid \mathbf{e})$

# Dealing with Evidence

- How can we use sampling to estimate  $P(\mathbf{y} \mid \mathbf{E} = \mathbf{e})$
- One approach is to sample from  $P(\mathbf{y} \mid \mathbf{E} = \mathbf{e})$  via *rejection sampling*
  - Draw  $\tilde{M}$  samples  $\mathbf{x} \sim P(\mathbf{X})$  as before
  - Discard (reject) samples that are inconsistent with  $\mathbf{E} = \mathbf{e}$
  - Resulting samples are from the posterior  $P(\mathbf{X} \mid \mathbf{e})$
- Expected number of particles that are kept is  $\tilde{M}P(\mathbf{e})$
- Requires  $\tilde{M} = M/P(\mathbf{e})$  samples to get  $M$  that are kept

# Dealing with Evidence

- How can we use sampling to estimate  $P(\mathbf{y} | \mathbf{E} = e)$
- One approach is to sample from  $P(\mathbf{y} | \mathbf{E} = e)$  via *rejection sampling*
  - Draw  $\tilde{M}$  samples  $\mathbf{x} \sim P(\mathbf{X})$  as before
  - Discard (reject) samples that are inconsistent with  $\mathbf{E} = e$
  - Resulting samples are from the posterior  $P(\mathbf{X} | e)$
- Expected number of particles that are kept is  $\tilde{M}P(e)$
- Requires  $\tilde{M} = M/P(e)$  samples to get  $M$  that are kept
- Unfortunately,  $P(e)$  is often small in practice (i.e., rare events)

# Dealing with Evidence

- How can we use sampling to estimate  $P(\mathbf{y} | \mathbf{E} = e)$
- One approach is to sample from  $P(\mathbf{y} | \mathbf{E} = e)$  via *rejection sampling*
  - Draw  $\tilde{M}$  samples  $\mathbf{x} \sim P(\mathbf{X})$  as before
  - Discard (reject) samples that are inconsistent with  $\mathbf{E} = e$
  - Resulting samples are from the posterior  $P(\mathbf{X} | e)$
- Expected number of particles that are kept is  $\tilde{M}P(e)$
- Requires  $\tilde{M} = M/P(e)$  samples to get  $M$  that are kept
- Unfortunately,  $P(e)$  is often small in practice (i.e., rare events)
- Alternatively, separately estimate  $P(\mathbf{y}, e)$  &  $P(e)$  and use Bayes' rule
  - Low relative error requires a number of samples that grows linearly with  $1/P(e)$
  - A bound on the absolute error for  $P(e)$ , which does not depend on  $P(e)$ , does not provide any bound on  $P(\mathbf{y}, e)/P(e)$

# Likelihood Weighting

- With forward sampling, evidence only affects sampling of descendents
- Forward sampling would reject particles that are inconsistent with evidence, but does not result in true posterior (e.g., 30% of the samples include  $(i^1, s^1)$ , which is the same as the prior over  $I$ )
- Basic idea: Weight each particle by the probability of the evidence variables taking the observed values (e.g., upweight  $(i^1, s^1)$  and downweight  $(i^0, s^1)$ )
- Now, we can estimate the conditional as

$$\hat{P}_{\mathcal{D}}(\mathbf{y} | \mathbf{e}) = \frac{\sum_{m=1}^M w^{(m)} \mathbb{1}(\mathbf{y}^{(m)} = \mathbf{y})}{\sum_{m=1}^M w^{(m)}}$$

where  $w^{(m)}$  is the product of probabilities of the evidence



# Likelihood Weighting

For  $m = 1, 2, \dots, M$ :

- ① Initialize  $w^{(m)} = 1$
- ② For each  $X_i$  in topological order:
  - If  $X_i$  is an evidence variable
    - $x_i^{(m)} \leftarrow e_i$
    - $w^{(m)} \leftarrow w^{(m)} \cdot P(x_i^{(m)} | \mathbf{u}_i^{(m)})$  where  $\mathbf{u}_i^{(m)}$  are sampled parents of  $x_i^{(m)}$
  - Else:
    - Sample  $x_i^{(m)}$  from  $P(x_i^{(m)} | \mathbf{u}_i^{(m)})$

# Likelihood Weighting

For  $m = 1, 2, \dots, M$ :

- ① Initialize  $w^{(m)} = 1$
- ② For each  $X_i$  in topological order:
  - If  $X_i$  is an evidence variable
    - $x_i^{(m)} \leftarrow e_i$
    - $w^{(m)} \leftarrow w^{(m)} \cdot P(x_i^{(m)} | \mathbf{u}_i^{(m)})$  where  $\mathbf{u}_i^{(m)}$  are sampled parents of  $x_i^{(m)}$
  - Else:
    - Sample  $x_i^{(m)}$  from  $P(x_i^{(m)} | \mathbf{u}_i^{(m)})$
- In the previous example, while 70% of the particles would involve  $(i^0, s^1)$ , their weight would only be 0.05
- Forward sampling is a special case of likelihood weighting with  $w^{(m)} = 1$

# Importance Sampling

- Often, we may not be able to sample from the *target distribution*  $P$ 
  - $P$  may not be known
  - Sampling from  $P$  may be computationally expensive (e.g., a posterior distribution for a Bayesian network)
- Instead, consider sampling from a (simpler) *proposal distribution*  $Q$ , s.t.  $Q > 0$  whenever  $P > 0$  (i.e., support of  $Q$  contains support of  $P$ )
- This gives rise to an equivalent expression for the expectation:

$$\mathbb{E}_P(f(\mathbf{X})) = \mathbb{E}_Q \left[ f(\mathbf{X}) \frac{P(\mathbf{X})}{Q(\mathbf{X})} \right]$$

# Importance Sampling

- Often, we may not be able to sample from the *target distribution*  $P$ 
  - $P$  may not be known
  - Sampling from  $P$  may be computationally expensive (e.g., a posterior distribution for a Bayesian network)
- Instead, consider sampling from a (simpler) *proposal distribution*  $Q$ , s.t.  $Q > 0$  whenever  $P > 0$  (i.e., support of  $Q$  contains support of  $P$ )
- This gives rise to an equivalent expression for the expectation:

$$\mathbb{E}_P(f(\mathbf{X})) = \mathbb{E}_Q \left[ f(\mathbf{X}) \frac{P(\mathbf{X})}{Q(\mathbf{X})} \right]$$

- Gives rise to the *unnormalized importance sampling* estimator: Given  $M$  samples from  $Q$ , weight each one by  $\frac{P(\mathbf{x}^{(m)})}{Q(\mathbf{x}^{(m)})}$ , yielding the estimate

$$\hat{\mathbb{E}}_{\mathcal{D}}[f] = \frac{1}{M} \sum_{m=1}^M f(\mathbf{x}^{(m)}) \frac{P(\mathbf{x}^{(m)})}{Q(\mathbf{x}^{(m)})}$$

# Unnormalized Importance Sampling

- (Admittedly, the name is confusing)
- Estimator is *unbiased*, i.e., for  $D$  sampled data sets

$$\mathbb{E}_D [\hat{\mathbb{E}}_D[f]] = \mathbb{E}_Q[f(\mathbf{X})w(\mathbf{X})] = \mathbb{E}_P[f(\mathbf{X})]$$

- The variance of the estimator becomes

$$\begin{aligned}\mathbb{V} &= \frac{1}{M} \mathbb{E}_Q[(f(\mathbf{X})w(\mathbf{X}))^2] - \mathbb{E}_Q[(f(\mathbf{X})w(\mathbf{X}))]^2 \\ &= \frac{1}{M} \mathbb{E}_Q[(w(\mathbf{X}))^2] - (\mathbb{E}_Q[(w(\mathbf{X}))])^2 \quad \text{when } f(\mathbf{X}) = 1\end{aligned}$$

(i.e., the variance of  $w(\mathbf{X}) = P(\mathbf{X})/Q(\mathbf{X})$ )

- Variance of the estimate decreases with
  - A larger number of samples
  - Choosing  $Q$  to be closer to  $P$ , specifically  $Q \propto |f(\mathbf{X})|P(\mathbf{X})$

# Normalized Importance Sampling

- However, we may not be able to evaluate  $P(\mathbf{x}^{(m)})$ 
  - May involve a posterior distribution over a Bayesian network
  - May require calculating a partition function (i.e., for an MRF)
- Suppose that we know  $\tilde{P}$ , i.e.,  $P$  up to a normalizing constant
- What if we use  $\tilde{P}$  in place of  $P$  when computing weights?
- Intuition: the weight is a random variable with expected value

$$\mathbb{E}_Q[w(\mathbf{X})] = \sum_{\mathbf{x}} Q(\mathbf{x}) \frac{\tilde{P}(\mathbf{x})}{Q(\mathbf{x})} = \sum_{\mathbf{x}} \tilde{P}(\mathbf{x}) = Z$$

# Normalized Importance Sampling

- We can write the objective as

$$\mathbb{E}_P[f(\mathbf{X})] = \frac{\mathbb{E}_Q[f(\mathbf{X})w(\mathbf{X})]}{\mathbb{E}_Q[w(\mathbf{X})]}$$

where

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

- Gives rise to *normalized importance sampling estimator*

$$\hat{\mathbb{E}}_{\mathcal{D}}[f] = \frac{\sum_{m=1}^M f(\mathbf{x}^{(m)})w(\mathbf{x}^{(m)})}{\sum_{m=1}^M w(\mathbf{x}^{(m)})} \quad w(\mathbf{x}^{(m)}) = \frac{\tilde{P}(\mathbf{x}^{(m)})}{Q(\mathbf{x}^{(m)})}$$

(*normalized* importance sampling uses the *unnormalized* distribution)

# Normalized Importance Sampling

- Not unbiased, but bias decreases as  $\frac{1}{M}$
- Typically, the variance is lower than that of the unnormalized estimator (but not always), outweighing added bias

$$\mathbb{V}_P[\hat{\mathbb{E}}_{\mathcal{D}}(f(\mathbf{X}))] \approx \frac{1}{M} \mathbb{V}_P[f(\mathbf{X})](1 + \mathbb{V}_Q[w(\mathbf{X})])$$

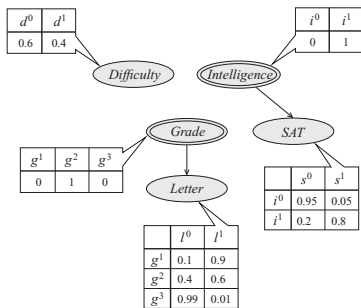
- We can use the above expression to measure the *effective sample size*, which is useful in deciding whether to keep sampling, i.e., for a set  $\mathcal{D}$  of  $M$  samples:

$$M_{\text{eff}} = \frac{M}{1 + \text{Var}[\mathcal{D}]}$$
$$\mathbb{V}[\mathcal{D}] = \sum_{m=1}^M w(\mathbf{x}^{(m)})^2 - \left( \sum_{m=1}^M w(\mathbf{x}^{(m)}) \right)^2$$



# Importance Sampling for Bayesian Networks

- Suppose that  $P$  is represented by a Bayesian network and that we are interested in event  $\mathbf{Z} = \mathbf{z}$
- The proposal distribution  $Q$  takes the form of a **mutilated network**  $\mathcal{B}_{\mathbf{Z}=\mathbf{z}}$  for instantiation  $\mathbf{Z} = \mathbf{z}$ , where:
  - Each  $Z_i \in \mathbf{Z}$  has no parents, and  $P(Z_i = z_i) = 1$  and zero otherwise
  - Parents and CPDs of all other variables  $X \notin \mathbf{Z}_i$  are unchanged



$$\mathcal{B}_{I=i^1, G=g^2}$$

# Ratio and Normalized Likelihood Weighting

- Determine  $P(\mathbf{Y} = \mathbf{y} \mid \mathbf{E} = \mathbf{e})$  by sampling from  $Q$  (mutilated network)
- Ratio likelihood weighting: Calculate  $P(\mathbf{Y} = \mathbf{y} \mid \mathbf{E} = \mathbf{e})$  using two runs of likelihood weighting, one for  $P(\mathbf{Y} = \mathbf{y}, \mathbf{E} = \mathbf{e})$  and the other for  $P(\mathbf{E} = \mathbf{e})$  (via Bayes' rule)

$$\hat{P}(\mathbf{y} \mid \mathbf{e}) = \frac{\hat{P}(\mathbf{y}, \mathbf{e})}{\hat{P}(\mathbf{e})} = \frac{1/M \sum_{m=1}^M w^{(m)}}{1/M' \sum_{m=1}^{M'} \bar{w}^{(m)}}$$

- Useful for single queries

# Ratio and Normalized Likelihood Weighting

- Determine  $P(\mathbf{Y} = \mathbf{y} \mid \mathbf{E} = \mathbf{e})$  by sampling from  $Q$  (mutilated network)
- Ratio likelihood weighting: Calculate  $P(\mathbf{Y} = \mathbf{y} \mid \mathbf{E} = \mathbf{e})$  using two runs of likelihood weighting, one for  $P(\mathbf{Y} = \mathbf{y}, \mathbf{E} = \mathbf{e})$  and the other for  $P(\mathbf{E} = \mathbf{e})$  (via Bayes' rule)

$$\hat{P}(\mathbf{y} \mid \mathbf{e}) = \frac{\hat{P}(\mathbf{y}, \mathbf{e})}{\hat{P}(\mathbf{e})} = \frac{1/M \sum_{m=1}^M w^{(m)}}{1/M' \sum_{m=1}^{M'} \bar{w}^{(m)}}$$

- Useful for single queries
- Normalized likelihood weighting: Calculate conditional probability for a *full distribution* over  $\mathbf{Y}$  with  $P(\mathbf{Y} \mid \mathbf{E} = \mathbf{e})$  as the target using normalized importance sampling evaluated at  $\tilde{P}(\mathbf{Y}, \mathbf{e})$  and the LW estimate

$$\hat{P}(\mathbf{y} \mid \mathbf{e}) = \frac{\sum_{m=1}^M \mathbb{1}(\mathbf{y}^{(m)} = \mathbf{y}) w^{(m)}}{\sum_{m=1}^M w^{(m)}}$$

# Performance of Importance Sampling

- Efficiency of importance sampling depends on how close proposal  $Q$  is to the target  $P$
- If the evidence is at the roots, then  $Q = P(\mathbf{Y} | e)$  (the posterior) and all samples have equal weight ( $P(e)$ ) (no need to compensate as evidence is sampled)
- If evidence is at the leaves, then  $Q = P(\mathbf{Y})$  (the prior), relying on the weights to correct for evidence
  - Samples are irrelevant if the prior and posterior differ significantly
  - Samples get small weight if evidence is unlikely

# Limitations of Monte Carlo

- In general, it is difficult to sample rare events, particularly in high-dimensions
- Importance sampling works poorly if  $Q$  is very different from  $P$
- Constructing  $Q$  to be similar to  $P$  can be difficult
  - A good proposal often requires knowledge of the analytic form of  $P$
- Instead . . . rather than use a fixed proposal  $Q$ , what if we could use an *adaptive* proposal?

- Intuition: “Fix” an initial sample by resampling some of the variables generated previously

# Markov Chain Monte Carlo

- Intuition: “Fix” an initial sample by resampling some of the variables generated previously
- General idea: generate an initial sample and then iteratively sample from each variable *given* previous samples of other variables
- Generates a *sequence* of samples

# Markov Chain Monte Carlo

- Intuition: “Fix” an initial sample by resampling some of the variables generated previously
- General idea: generate an initial sample and then iteratively sample from each variable *given* previous samples of other variables
- Generates a *sequence* of samples
- In the beginning, likely not sampling from  $P$
- Over time, sampling distribution gets closer to  $P$



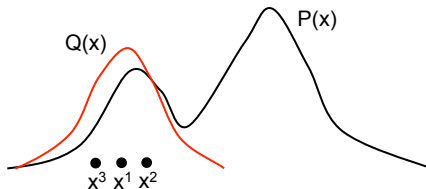
# Markov Chain Monte Carlo

- Intuition: “Fix” an initial sample by resampling some of the variables generated previously
- General idea: generate an initial sample and then iteratively sample from each variable *given* previous samples of other variables
- Generates a *sequence* of samples
- In the beginning, likely not sampling from  $P$
- Over time, sampling distribution gets closer to  $P$
- Useful in both directed and undirected models (unlike likelihood weighting and other forward sampling methods)

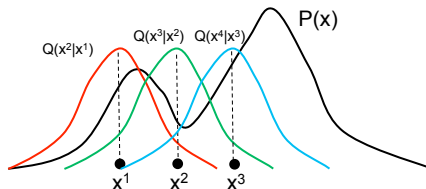
# Markov Chain Monte Carlo

- MCMC algorithms employ adaptive proposals
  - Instead of  $Q(x)$ , MCMC uses  $Q(x' | x)$  where  $x$  is the previous sample
  - As  $x$  changes, so does  $Q(x' | x)$  (as a function of  $x'$ )

Importance sampling with  
a (bad) proposal  $Q(x)$



MCMC with adaptive  
proposal  $Q(x'|x)$



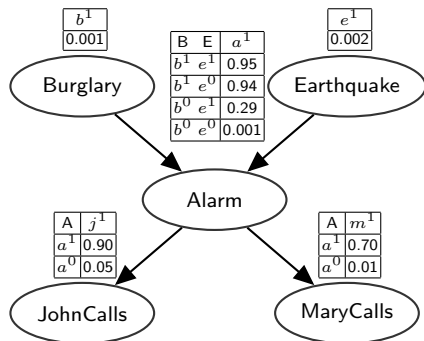
# Gibbs Sampling

Input: Random variables  $\mathbf{X}$ , factors  $\Phi$ , initial state distribution  $P^{(0)}(\mathbf{X})$ , and number of time steps  $T$

- ① Sample  $\mathbf{x}^{(0)} \sim P^{(0)}(\mathbf{X})$
- ② for  $t = 1, \dots, T$ 
  - $\mathbf{x}^{(t)} \leftarrow \mathbf{x}^{(t-1)}$
  - For each  $X_i \in \mathbf{X}$ : Sample  $x_i^{(t)} \sim P(X_i | \mathbf{x}_{-i})$
- ③ Return  $\mathbf{x}^{(0)}, \dots, \mathbf{x}^{(T)}$

- During each round, we sample one random variable conditioned on other variables fixed to the most recent sampled value
  - This only involves a subset of the factors
  - We only need to look at some values of those factors
- We can consider evidence by first reducing the factors
- Appeal: Unlike forward sampling, each variable is influenced by the others, including downstream evidence

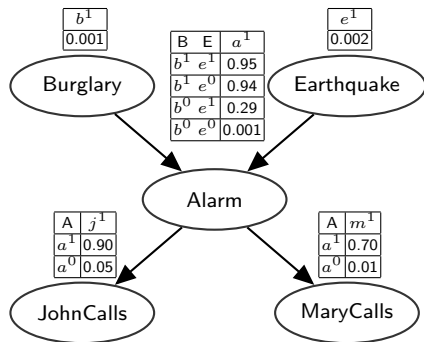
# Gibbs Sampling: Example



t	$B$	$E$	$A$	$J$	$M$
0	F	F	F	F	F
1					
2					
3					
4					

- Consider the alarm network example
- Assume that we sample variables in the order  $B, E, A, J, M$
- Initialize all variables at  $t = 0$  to False

# Gibbs Sampling: Example



t	B	E	A	J	M
0	F	F	F	F	F
1	F				
2					
3					
4					

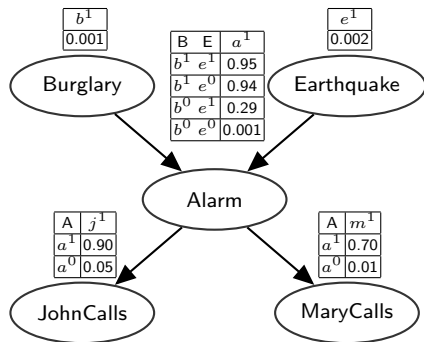
- Sample  $P(B | A, E)$  at  $t = 1$  via Bayes rule

$$P(B | A, E) \propto P(A | B, E)P(B)$$

- $A = \text{False}, E = \text{False}$ :

- $P(B = \text{T} | A = \text{F}, E = \text{F}) \propto 0.06 \cdot 0.01 = 0.0006$
- $P(B = \text{F} | A = \text{F}, E = \text{F}) \propto 0.999 \cdot 0.999 = 0.9980$

# Gibbs Sampling: Example



t	B	E	A	J	M
0	F	F	F	F	F
1	F	T			
2					
3					
4					

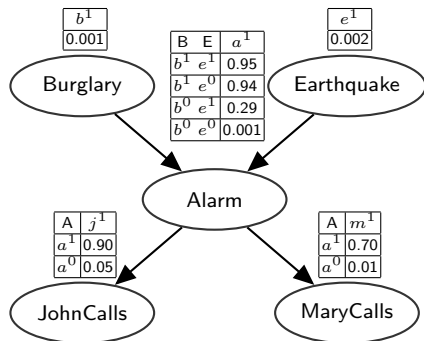
- Sample  $P(E | A, B)$  at  $t = 1$  via Bayes rule

$$P(E | A, B) \propto P(A | B, E)P(E)$$

- $A = \text{False}$ ,  $B = \text{False}$ :

- $P(E = T | A = F, B = F) \propto 0.71 \cdot 0.02 = 0.0142$
- $P(E = F | A = F, B = F) \propto 0.999 \cdot 0.998 = 0.9970$

# Gibbs Sampling: Example



t	B	E	A	J	M
0	F	F	F	F	F
1	F	T	F		
2					
3					
4					

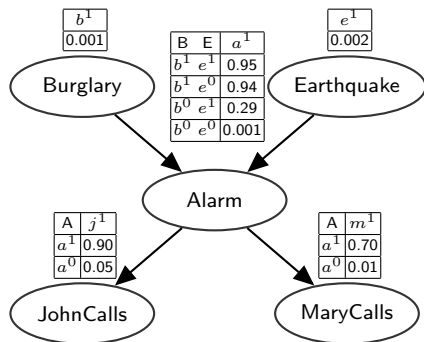
- Sample  $P(A | B, E, J, M)$  at  $t = 1$  via Bayes rule

$$P(A | B, E, J, M) \propto P(J | A)P(M | A)P(A | B, E)$$

- $(B, E, J, M) = (F, T, F, F)$ :
  - $P(A = T | B = F, E = T, J = F, M = F) \propto 0.1 \cdot 0.3 \cdot 0.29 = 0.0087$
  - $P(A = F | B = F, E = T, J = F, M = F) \propto 0.95 \cdot 0.99 \cdot 0.71 = 0.6678$



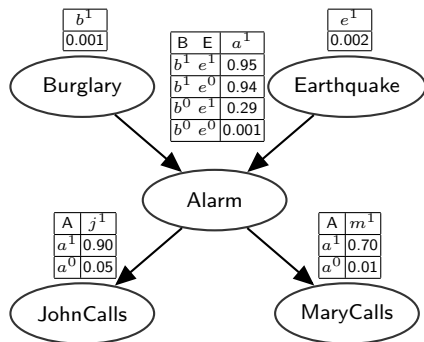
# Gibbs Sampling: Example



t	B	E	A	J	M
0	F	F	F	F	F
1	F	T	F	T	
2					
3					
4					

- Sample  $P(J | A)$  at  $t = 1$  directly (no need for Bayes rule)
- $A = F$ :
  - $P(J = T | A = F) = 0.05$
  - $P(J = F | A = F) = 0.95$

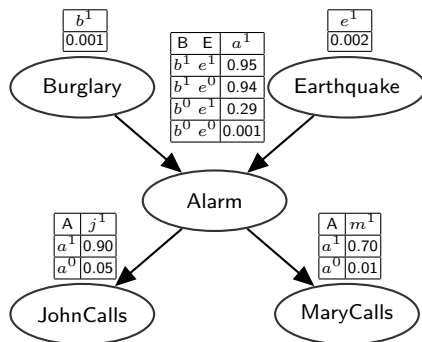
# Gibbs Sampling: Example



t	B	E	A	J	M
0	F	F	F	F	F
1	F	T	F	T	F
2					
3					
4					

- Sample  $P(M | A)$  at  $t = 1$  directly (no need for Bayes rule)
- $A = F$ :
  - $P(M = T | A = F) = 0.01$
  - $P(M = F | A = F) = 0.99$

# Gibbs Sampling: Example

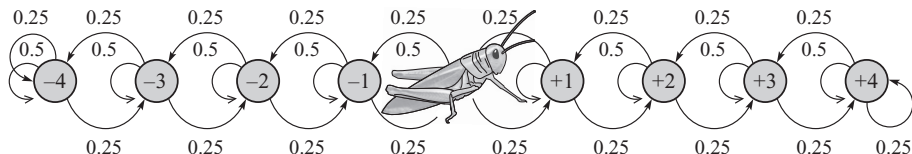


t	B	E	A	J	M
0	F	F	F	F	F
1	F	T	F	T	F
2	F	T	T	T	T
3	T	F	T	F	T
4	T	F	T	F	F

- For  $t = 2, 3, \dots$ , we repeat this process

# Markov Chains

- A **Markov chain** is defined via a *state space*  $\text{Val}(\mathbf{X})$  and a *transition model*  $\mathcal{T}(x \rightarrow x')$  specifying probability of going from  $x$  to  $x'$
- Defined in terms of a graph of states
  - This graph is different from the original graphical model
  - Nodes are possible assignments to  $\mathbf{X}$
- Sampler takes a *random walk* through the states



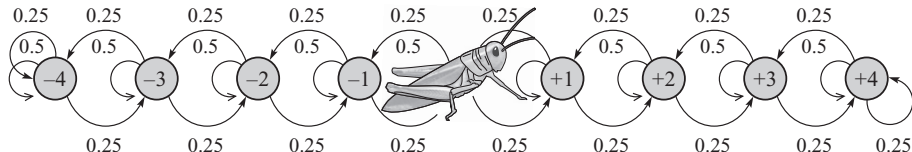
**Figure:** Markov chain where states are  $x \in \{-4, -3, \dots, 3, 4\}$ .

# Markov Chains

- Random samples dictated by chain dynamics

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

- Asymptotic behavior: What does the long-term distribution  $P^{(t)}$  for  $t \gg 1$  look like? (Grasshopper example converges to uniform)
- Asymptotic behavior: We want  $P^{(t)}$  to converge to  $P$



**Figure:** Markov chain where states are  $x \in \{-4, -3, \dots, 3, 4\}$ .

# Markov Chains: Asymptotic Behavior

- Distribution  $\pi(\mathbf{X})$  is a **stationary distribution** for Markov chain  $\mathcal{T}$  if

$$\pi(\mathbf{X} = \mathbf{x}') = \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \pi(\mathbf{X} = \mathbf{x}) \mathcal{T}(\mathbf{x} \rightarrow \mathbf{x}')$$

(intuitively,  $P^{(t+1)}$  is close to  $P^{(t)}$ )

- If  $A_{i,j} = \mathcal{T}(\mathbf{x} \rightarrow \mathbf{x}')$ , then the stationary distribution is an eigenvector of  $A$  with eigenvalue 1

# Markov Chains: Asymptotic Behavior

- Distribution  $\pi(\mathbf{X})$  is a **stationary distribution** for Markov chain  $\mathcal{T}$  if

$$\pi(\mathbf{X} = \mathbf{x}') = \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \pi(\mathbf{X} = \mathbf{x}) \mathcal{T}(\mathbf{x} \rightarrow \mathbf{x}')$$

(intuitively,  $P^{(t+1)}$  is close to  $P^{(t)}$ )

- If  $A_{i,j} = \mathcal{T}(\mathbf{x} \rightarrow \mathbf{x}')$ , then the stationary distribution is an eigenvector of  $A$  with eigenvalue 1
- A *periodic Markov chain* exhibits cyclic behavior
- A *reducible Markov chain* Has a non-unique stationary distribution that depends on  $P^{(0)}$  (i.e., there are sets of states that can not be reached from one another)

# Markov Chains: Asymptotic Behavior

- Distribution  $\pi(\mathbf{X})$  is a **stationary distribution** for Markov chain  $\mathcal{T}$  if

$$\pi(\mathbf{X} = \mathbf{x}') = \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \pi(\mathbf{X} = \mathbf{x}) \mathcal{T}(\mathbf{x} \rightarrow \mathbf{x}')$$

- We want a Markov chain with a unique stationary distribution
- A Markov chain is *regular* (*ergodic*) if  $\exists k$  s.t.  $\forall \mathbf{x}, \mathbf{x}' \in \text{Val}(\mathbf{X})$ , the probability of getting from  $\mathbf{x}$  to  $\mathbf{x}'$  in exactly  $k$  steps is  $> 0$ 
  - Sufficient: (i) possible to get from any state to any other state via positive probability path; and (ii) self-loops
- A Markov chain that is regular has a unique stationary distribution



# Markov Chains: Asymptotic Behavior

A stationary distribution  $\pi$  satisfies

$$\pi(x^1) = 0.25\pi(x^1) + 0.5\pi(x^3)$$

$$\pi(x^2) = 0.7\pi(x^2) + 0.5\pi(x^3)$$

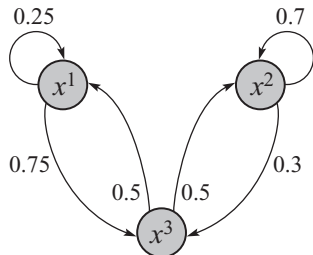
$$\pi(x^3) = 0.75\pi(x^1) + 0.3\pi(x^2)$$

and

$$\pi(x^1) + \pi(x^2) + \pi(x^3) = 1,$$

yielding

$$\pi(x^1) = 0.2 \quad \pi(x^2) = 0.5 \quad \pi(x^3) = 0.3$$



# Markov Chains Monte Carlo (MCMC) Sampling

Markov chain dynamics give rise to the MCMC sampling process

---

## Algorithm 12.5 Generating a Markov chain trajectory

---

**Procedure** MCMC-Sample (

$P^{(0)}(\mathbf{X})$ , // Initial state distribution

$\mathcal{T}$ , // Markov chain transition model

$T$  // Number of time steps

)

1 Sample  $\mathbf{x}^{(0)}$  from  $P^{(0)}(\mathbf{X})$

2 **for**  $t = 1, \dots, T$

3     Sample  $\mathbf{x}^{(t)}$  from  $\mathcal{T}(\mathbf{x}^{(t-1)} \rightarrow \mathbf{X})$

4     **return**  $\mathbf{x}^{(0)}, \dots, \mathbf{x}^{(T)}$

---

# Markov Chains: Transition Models

- The state space for graphical models has a factorized structure (each state is an assignment to several variables)
- Rather than transitions between general states, we can consider transitions that update only a single state element at a time
- Gives rise to set of transition models  $\{\mathcal{T}_1, \dots, \mathcal{T}_k\}$  (kernels)
- A *multi-kernel Markov chain* cycles through these kernels
- Graphical model: One  $\mathcal{T}_i((\mathbf{x}_{-i}, x_i) \rightarrow (\mathbf{x}_{-i}, x'_i))$  for each  $X_i$
- While individual kernels may not be regular, the combination may be

# Gibbs Sampling and Markov Chains

- A *Gibbs chain* is defined by the following kernel

$$\mathcal{T}_i((\mathbf{x}_{-i}, x_i) \rightarrow (\mathbf{x}_{-i}, x'_i)) = P(x'_i | \mathbf{x}_{-i})$$

- Each step of Gibbs sampling visits a new state in the chain
- Factors allow efficient calculations at each step
- Over time, sampler converges to the stationary distribution, i.e.,  $\pi(\mathbf{X}) = P(\mathbf{X} | \mathbf{e})$ , where  $\mathbf{X} = \mathcal{X} - \mathbf{E}$

- It may take MCMC a long time to *mix*, i.e., converge to the stationary distribution
- Often, we throw out early samples that correspond to the “burn-in” period