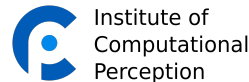


PROBABILISTIC MODELS – PART 4b: APPROXIMATE INFERENCE IN BAYESIAN NETWORKS

Gerhard Widmer

Institute of Computational Perception
Johannes Kepler University
Linz, Austria

gerhard.widmer@jku.at
www.cp.jku.at/people/widmer



October 14, 2025

Presentation partly based on and inspired by [Koller & Friedman, 2009] and [Russell & Norvig, 2021], including the use of some figures from their books and/or lecture slides.

Many thanks to Daphne Koller, Nir Friedman, Stuart Russell, and Peter Norvig for making these available (`pgm.stanford.edu`; `aima.cs.berkeley.edu`).

Do not distribute!

Goals of this Lecture

- ▶ Discuss the need for approximate reasoning in BNs
- ▶ Introduce a specific class of approximate reasoning algorithms: Stochastic Sampling
- ▶ Show how to sample events from a network
- ▶ Show how to answer probabilistic queries based on such samples
- ▶ Introduce the general concept of Importance Sampling
- ▶ Introduce the framework of Markov Chain Monte Carlo (MCMC) methods
- ▶ Describe Gibbs Sampling as an instance of MCMC
- ▶ Briefly characterise a general strategy for constructing Markov chains for arbitrary target distributions: the Metropolis-Hastings Algorithm

Outline

1 Motivation

2 Simple Sampling Strategies

Drawing Samples from a Bayesian Network: Forward Sampling

Estimating the Full Joint Distribution

Answering Probabilistic Queries:

- Rejection Sampling
- Likelihood Weighting

3 The General Importance Sampling Framework

The General Idea

Likelihood Weighting as an Instance of Importance Sampling

4 Markov Chain Monte Carlo (MCMC) Sampling

The General Idea

Gibbs Sampling

Gibbs Sampling as an Instance of MCMC

Constructing Markov Chains: The Metropolis-Hastings Algorithm

Motivation

Remember: The Bad News

- ▶ It can be shown that the problem of inference in graphical models is \mathcal{NP} -hard, and therefore requires **exponential time**, in the worst case.

The (Partly) Good News:

- ▶ For some sparsely connected networks (i.e., networks where variables tend to have few parents), inference can be done quite efficiently with optimised algorithms (e.g., 📖 *Variable Elimination*)

But: For many real-world problems, exact inference is impossible.

One Possible Solution: Approximate Inference

General Idea:

- ▶ Look for ways to compute **approximate** answer with tractable computational effort.

Different Approaches:

- ▶ **Optimisation-based Methods ('Variational Inference')**: Do *exact* inference in a distribution Q that is 'simpler' than target distribution P ; search for best approximation distribution Q in a class \mathcal{Q} of 'easy' distributions.
- ▶ **Stochastic Sampling Methods**: Do not perform exact inference; answer queries by generating random sample events according to P and estimating probabilities via counting.

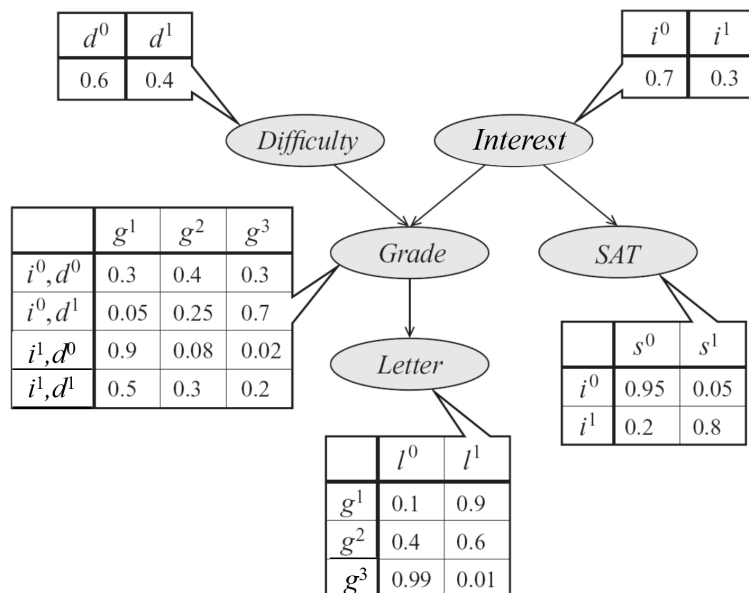
In this Class:

- ▶ **Stochastic Sampling** (also known as '**Particle-Based Methods**')

Approximate Inference via Stochastic Sampling

Basic Idea:

- ▶ Instead of exactly computing a probability distribution $P(\mathcal{X})$, approximate it via a limited set of points randomly drawn from P (if that can be done)
- ▶ These points are called **(random) samples** or **particles**
- ▶ Estimate any desired probability $P(x)$ from these samples.



⇒ sampled particles
(atomic events):

$$\begin{pmatrix} d^1 & i^0 & g^2 & s^1 & l^0 \\ d^1 & i^1 & g^1 & s^0 & l^1 \\ d^0 & i^0 & g^3 & s^0 & l^1 \\ d^1 & i^0 & g^1 & s^1 & l^1 \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \end{pmatrix}$$

Sampling Events from a Bayesian Network

Question:

- ▶ How to generate a random **atomic event** according to a distribution $P_{\mathcal{B}}(\mathcal{X})$ represented as a Bayesian Network \mathcal{B} ?

The **Forward Sampling Algorithm**:

- ▶ Sample variables independently, one by one, starting at the ‘top’

The Forward Sampling Algorithm

Given:

- ▶ A Bayesian Network \mathcal{B} over random variables \mathcal{X}

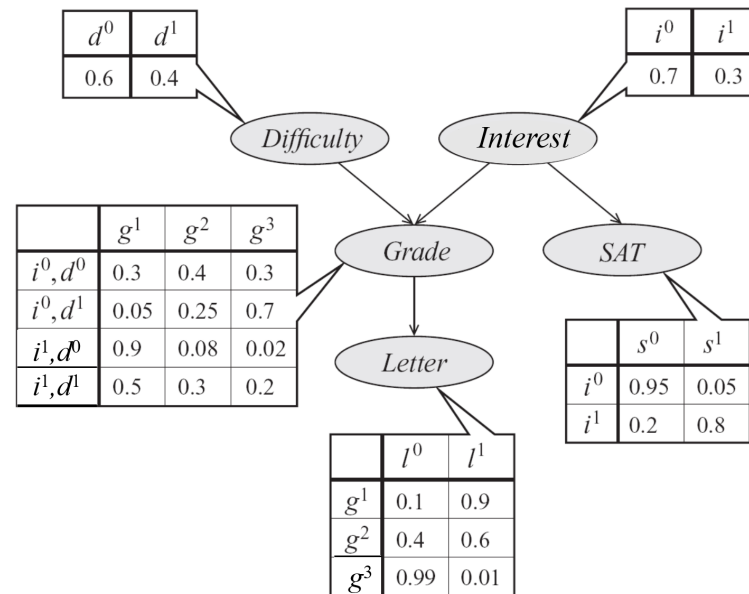
Compute:

- ▶ A random event x , drawn according to the probability distribution $P_{\mathcal{B}}(\mathcal{X})$

Algorithm:

- ▶ Choose a **topological ordering** $\mathbf{X} = \{X_1, X_2, \dots, X_N\}$ of the \mathcal{X}
- ▶ For $i = 1$ to N :
Sample a value x_i from $Val(X_i)$ **according to the conditional distribution** $P(X_i \mid \text{pa}(X_i))$, given the parent values $\text{pa}(X_i)$ sampled previously.

Forward Sampling: An Example



Result:

Random sample (event)

$(d^1, i^0, g^3, s^0, l^1)$

Topological ordering: D, I, G, S, L

Sampling steps:

- 1 Sample D : randomly choose a value from $Val(D) = \{d^0, d^1\}$ acc. to $P(D) = [0.6; 0.4]$
 \Rightarrow suppose we get $\underline{d^1}$
- 2 Sample I from $\{i^0, i^1\}$ acc. to $P(I) = [0.7; 0.3]$
 \Rightarrow suppose we get $\underline{i^0}$
- 3 Sample G from $\{g^1, g^2, g^3\}$ acc. to $P(G \mid \underline{d^1}, \underline{i^0}) = [0.05; 0.25; 0.7]$
 \Rightarrow suppose we get $\underline{g^3}$
- 4 Sample S from $\{s^0, s^1\}$ acc. to $P(S \mid \underline{i^0}) = [0.95; 0.05]$
 \Rightarrow suppose we get $\underline{s^0}$
- 5 Sample L from $\{l^0, l^1\}$ acc. to $P(L \mid \underline{g^3}) = [0.99; 0.01]$
 \Rightarrow suppose we get $\underline{l^1}$

Proof: Forward Sampling is Consistent

To be shown:

- ▶ Forward Sampling is **consistent**
- ▶ **In words:** The probability distribution P_S according to which the samples are generated (the **Sampling Distribution**) is equal to the true distribution P_B represented by the Bayesian network.

Proof (trivial):

- ▶ Let $P_S(x_1, \dots, x_n)$ = probability that a specific event (x_1, \dots, x_n) is generated by the Forward Sampling algorithm
- ▶ Observation 1: Each variable X_i is sampled according to $P(X_i \mid \text{pa}(X_i))$
- ▶ Observation 2: Each variable is sampled *independently* of the other variables (except for using the variable's previously sampled parent values)
- ▶ Consequence:

$$P_S(x_1, \dots, x_n) = \prod_{i=1}^n P(X_i \mid \text{pa}(X_i)) = P_B(x_1, \dots, x_n). \quad \square$$

Computational Complexity of Forward Sampling

Let

$$n = |\mathcal{X}|$$

$$k = \max_i |\text{Pa}(X_i)|$$

$$d = \max_i |\text{Val}(X_i)|$$

Then the cost for

- ▶ indexing into the appropriate CPD $P(X_i \mid \text{pa}(X_i))$ of a variable X_i , given specific values of its parents, is $O(|\text{Pa}(X_i)|) \leq O(k)$
- ▶ sampling a value for a discrete variable X_i from a CPD $P(X_i \mid \text{pa}(X_i))$ is $O(\log |\text{Val}(X_i)|) \leq O(\log d)$ (via binary search)
- ▶ sampling a complete event is $O(nk \log d)$

👉 **Generating one random sample event can be done in time**

$$O(nk \log d)$$

Approximate Inference (1): Estimating the Full Joint Distribution

Problem:

- ▶ Estimate the probability $P_{\mathcal{B}}(x_1, \dots, x_n)$ of an **atomic event** (x_1, \dots, x_n) , where $\{X_1, \dots, X_n\} = \mathcal{X}$

Solution:

- ▶ Sample a set \mathcal{D} of M events randomly from \mathcal{B} via Forward Sampling
- ▶ Count the number of times K that our event (x_1, \dots, x_n) appears in \mathcal{D}
- ▶ Estimate $P_{\mathcal{B}}(x_1, \dots, x_n)$ via its relative frequency:

$$\hat{P}_{\mathcal{D}}(x_1, \dots, x_n) = \frac{K}{M}$$

👉 Can estimate the entire Full Joint Distribution (theoretically ...)

Approximate Inference (2): Estimating the Probability of an Arbitrary Event

Trivial generalisation to probability estimation for **general (not necessarily atomic)** events:

Problem:

- ▶ Estimate the probability $P_{\mathcal{B}}(x_1, \dots, x_k)$ of an event (x_1, \dots, x_k) , where $\{X_1, \dots, X_k\} \subseteq \mathcal{X}$

Solution:

- ▶ Sample a set \mathcal{D} of M atomic events randomly from \mathcal{B} via Forward Sampling
- ▶ Count the number K of events e_i in \mathcal{D} that **match** our event (x_1, \dots, x_k) (that is: where $e_i \langle X_1, \dots, X_k \rangle = (x_1, \dots, x_k)$)
- ▶ Estimate $P_{\mathcal{B}}(x_1, \dots, x_k)$ via the relative frequency

$$\hat{P}_{\mathcal{D}}(x_1, \dots, x_k) = \frac{K}{M}$$

Examples

Query 1 (atomic event):

$$P(d^1, i^0, g^3, s^0, l^1) = ?$$

Suppose we sample $M = 1000$ events, of which $K = 27$ events have precisely the values $(d^1, i^0, g^3, s^0, l^1)$

$$\hat{P}(d^1, i^0, g^3, s^0, l^1) = \frac{27}{1000} = 0.027$$

Query 2:

$$P(i^0, l^1) = ?$$

Suppose we sample $M = 1000$ events, of which $K = 199$ events have values $\{i^0, l^1\}$ for $\{I, L\}$ (and arbitrary values for the other variables $\{D, G, S\}$)

$$\hat{P}(i^0, l^1) = \frac{199}{1000} = 0.199$$

Approximate Inference (3): Probabilistic Queries

Remember:

A **Probabilistic Query** involves computing the conditional distribution

$$P(\mathbf{X} \mid \mathbf{E} = e) \quad (\text{or } P(\mathbf{X} \mid e) \text{ for short})$$

for some sets of variables $\mathbf{X}, \mathbf{E} \subseteq \mathcal{X}$ and a specific value assignment $\mathbf{E} = e$.

👉 Estimating $P(\mathbf{X} \mid e)$ means estimating $P(x \mid e)$ for all value combinations x .

Overview of the methods to be presented in the following:

- ▶ Simple Idea 1: Rejection Sampling (RS)
- ▶ Simple Idea 2: Likelihood Weighting (LW)
- ▶ More general view of LW: Importance Sampling
- ▶ Less simple Idea 3: Getting closer to the posterior via MCMC
- ▶ Specific realisation: Gibbs Sampling
- ▶ General method to create appropriate Markov Chains:
The Metropolis-Hastings Algorithm

Estimating $P(X|e)$ [1]: Rejection Sampling

Idea:

- ▶ Produce a set of samples from $P_{\mathcal{B}}$
- ▶ Keep only samples consistent with evidence $E = e$
- ▶ Estimate probabilities by counting relative frequencies (for each x)

Justification:

$$P(x | e) = \frac{P(x, e)}{P(e)}$$

and we can estimate both $P(x, e)$ and $P(e)$ as before:

$$\hat{P}(x | e) = \frac{\frac{N(x, e)}{M}}{\frac{N(e)}{M}} = \frac{N(x, e)}{N(e)}$$

where

$N(x, e)$... number of samples with values (x, e)
$N(e)$... number of samples with values e
M	... total number of samples.

Estimating $P(\mathbf{X} | e)$ [1]: Rejection Sampling

The Rejection Sampling Algorithm

Given:

- ▶ A Bayesian network \mathcal{B} over \mathcal{X}
- ▶ A query $P(\mathbf{X} | e)$ with $\mathbf{X}, \mathbf{E} \subseteq \mathcal{X}$

Algorithm:

- 1 **Sample:** Probabilistically sample a set \mathcal{D} of M events from \mathcal{B}
- 2 **Filter / Reject:**
Discard (“reject”) all samples not compatible with the evidence e
 \Rightarrow the remaining samples $\mathcal{D}_e \subseteq \mathcal{D}$ represent the distribution $P(\mathbf{X} | e)$
- 3 **Estimate:** For each value combination $\mathbf{X} = \mathbf{x}$, count the number K of samples $\in \mathcal{D}_e$ that match \mathbf{x} . Estimate $P(\mathbf{x} | e)$ as

$$\hat{P}(\mathbf{x} | e) = \frac{K}{|\mathcal{D}_e|}$$

Rejection Sampling: An Example

Query: $P(i^1 \mid g^2, s^1, l^0) = ?$

Suppose

- ▶ we sample a set \mathcal{D} of $M = 1000$ events
- ▶ 916 of these are *not* consistent with evidence $G = g^2 \wedge S = s^1 \wedge L = l^0$
 $\Rightarrow \mathcal{D}_e$ contains 84 events
- ▶ $K = 31$ of these have $I = i^1$

Answer:

$$\hat{P}(i^1 \mid g^2, s^1, l^0) = \frac{31}{84} \approx 0.369$$

Rejection Sampling: Pros and Cons

Advantages:

- ▶ Extremely simple algorithm
- ▶ Is consistent: produces provably correct estimate (in the limit)

Disadvantages / Problems:

- ▶ Wasteful: Immediately discards a large part of the generated samples
- ▶ More precisely: Expected number of samples that are *not* rejected (in sample set \mathcal{D} of size M) is $M \cdot P(e)$
- ▶ Example: with $P(e) = 0.001$, would need to generate 100,000 samples to obtain 100 samples useful for estimation!
- ▶ Probability of evidence e decreases exponentially with $|E|$

Source of the problem:

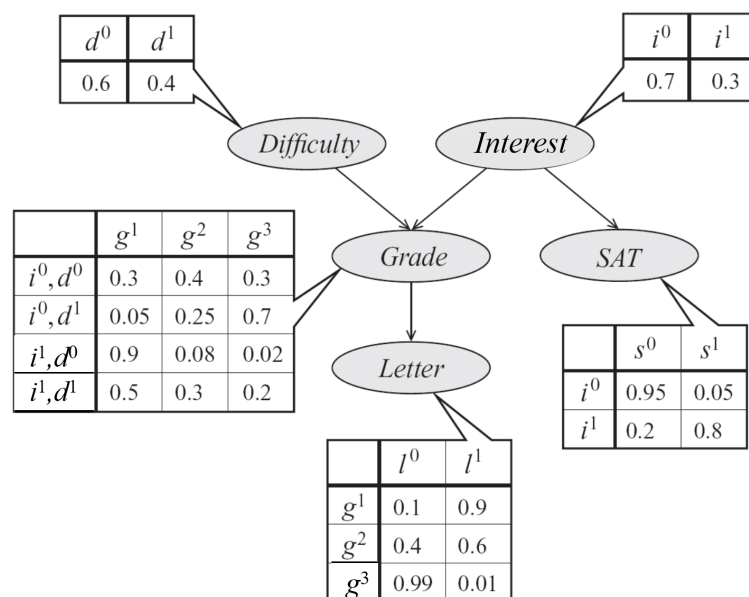
- ▶ Want to estimate the *posterior distribution* $P(\mathbf{X} \mid e)$
- ▶ ... but draw our samples \mathcal{D} from the *full joint* distribution $P(\mathcal{X})$ (may be very different from the posterior)
- ▶ Need to fix this by reducing \mathcal{D} to \mathcal{D}_e , which represents our target $P(\mathbf{X} \mid e)$



In the following: Try to sample from distributions closer to the posterior ...

Estimating $P(X|e)$ [2]: Likelihood Weighting

Motivation:



Query: $P(L | d^1, s^1) = ?$

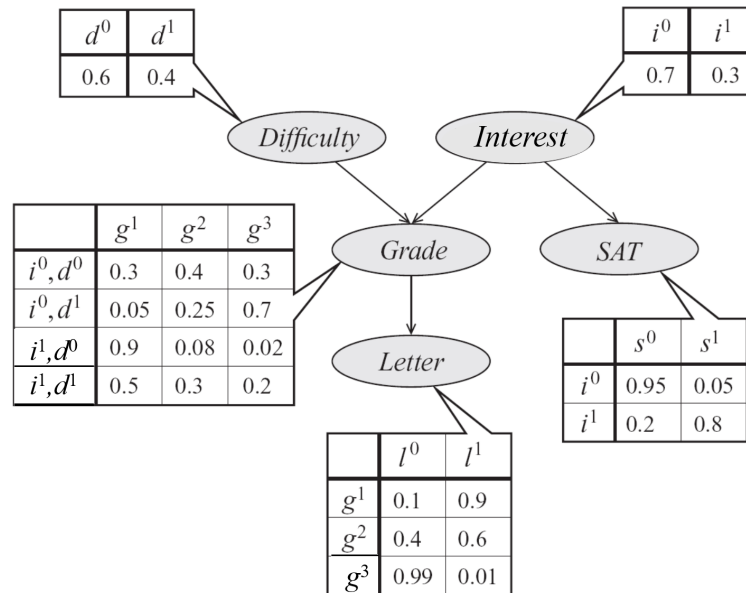
- ▶ Forward Sampling might start by generating a value d^0 for D
- ▶ No matter how the sampling proceeds, the resulting event will eventually be rejected!



... so why not just:

- ▶ Set evidence variables E to their observed values e
and only sample the remaining variables $\mathcal{X} - E$ to directly obtain \mathcal{D}_e

Problem with Naive Approach



Query: $P(L \mid d^1, s^1) = ?$

Naive Idea:

- 1 Set $D = d^1, S = s^1$
- 2 Sample I from its prior distribution
- 3 Sample G and L accordingly.

- ▶ All samples will have $D = d^1, S = s^1$, as desired
- ▶ Around 30% of the samples will have $I = i^1$
(because I is sampled from its prior distribution)
- ▶ **Problem:** True probability of $I = i^1$ in situations with $S = s^1$ is much higher:

$$P(i^1 \mid s^1) \gg P(i^1)$$

- ▶ Naive algorithm is unaware of this: it ignores evidence s^1 when sampling I
- ▶ Sampling distribution will be wrong and not reflect $P(I, G, L \mid d^1, s^1)$

Likelihood Weighting (LW): The Intuition

Consider this:

- ▶ Run standard Forward Sampling many times
- ▶ $P(s^1 | i^1) = .8$
 \Rightarrow Samples starting with $I = i^1$ would generate $S = s^1$ in 80% of the cases
- ▶ $P(s^1 | i^0) = .05 \Rightarrow$ Samples with $I = i^0$ would generate $S = s^1$ in only 5%
- ▶ Observing $S = s^1$ (evidence) in $I = i^1$ is $16\times$ as probable as when $I = i^0$

Idea:

- ▶ Simulate this by **sample weights** that reflect the probability of the evidence in a given context
- ▶ Example: Sample with $I = i^1$ and fixed evidence $S = s^1$ should be worth 80% of a sample, one with $I = i^0$ and $S = s^1$ only 5%
- ▶ Weight contributed by evidence variable S corresponds to $P(s^1 | I)$
 (in general: $P(e_i | \text{pa}(E_i))$)
- ▶ Variables are sampled independently \Rightarrow weight for entire sample = product of individual weights contributed by evidence variables.

Weighted Particles

Likelihood Weighting (LW): The Algorithm

The Likelihood Weighting Algorithm

Given:

- ▶ A Bayesian network \mathcal{B} over \mathcal{X}
- ▶ A query $P(\mathbf{X} \mid \mathbf{e})$ with $\mathbf{X}, \mathbf{E} \subseteq \mathcal{X}$

Algorithm:

- 1 **Sample:** Generate a set \mathcal{D}_e of **weighted samples** from \mathcal{B} (all consistent with e), via algorithm **SAMPLE-LW** (see next slide)
- 2 **Estimate:** For each value combination $\mathbf{X} = \mathbf{x}$, select the subset $\mathcal{D}_e^{\mathbf{x}} \subseteq \mathcal{D}_e$ of samples that match \mathbf{x} . Estimate $P(\mathbf{x} \mid \mathbf{e})$ as

$$\hat{P}(\mathbf{x} \mid \mathbf{e}) = \frac{\sum_{s_i \in \mathcal{D}_e^{\mathbf{x}}} w[s_i]}{\sum_{s_i \in \mathcal{D}_e} w[s_i]}$$

where $w[s_i]$ is the weight of sample s_i

Drawing a Weighted Sample

The SAMPLE-LW Algorithm

Given:

- ▶ A Bayesian network \mathcal{B} over \mathcal{X}
- ▶ Fixed evidence $\mathbf{E} = e$

Compute:

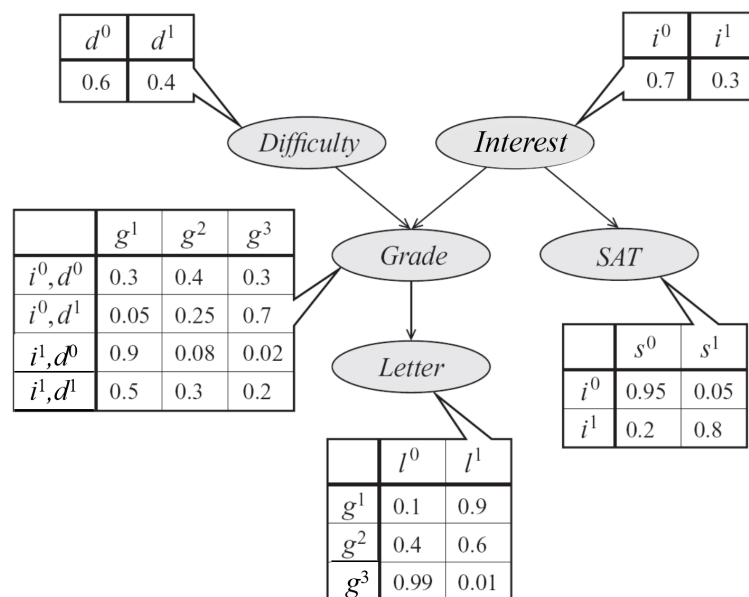
- ▶ A weighted sample $\langle (\mathbf{y}, e), w \rangle$ (where $\mathbf{Y} \cup \mathbf{E} = \mathcal{X}$)

Algorithm:

- ▶ Choose a **topological ordering** $\mathbf{X} = \{X_1, X_2, \dots, X_N\}$ of the \mathcal{X}
- ▶ Initialise weight w of current sample to $w = 1.0$
- ▶ For $i = 1$ to N :
 - if** X_i is one of the evidence variables $E_j \in \mathbf{E}$:
 - set X_i to its given evidence value e_j
 - update weight: $w \leftarrow w \times P(e_j \mid \text{pa}(X_i))$
 - else**
 - sample a value for X_i from $\text{Val}(X_i)$ according to $P(X_i \mid \text{pa}(X_i))$
 - (as in Forward Sampling)

SAMPLE-LW: An Example

Query: $P(I \mid g^3, s^1) = ?$



Result:

Weighted random sample:

$\langle (d^1, i^0, g^3, s^1, l^1), 0.035 \rangle$

Topological ordering: D, I, G, S, L

- 1 Initialise weight: $w \leftarrow 1.0$
- 2 Sample D from $\{d^0, d^1\}$ acc. to $P(D) = [0.6; 0.4]$
 \Rightarrow suppose we get $\underline{d^1}$
- 3 Sample I from $\{i^0, i^1\}$ acc. to $P(I) = [0.7; 0.3]$
 \Rightarrow suppose we get $\underline{i^0}$
- 4 G is given (evidence): Set G to $\underline{g^3}$ and update weight: $w \leftarrow w \times P(g^3 | d^1, i^0) = 1.0 \times 0.7 = 0.7$
- 5 S is given (evidence): Set S to $\underline{s^1}$ and update weight: $w \leftarrow w \times P(s^1 | i^0) = 0.7 \times 0.05 = 0.035$
- 6 Sample L from $\{l^0, l^1\}$ acc. to $P(L | \underline{g^3}) = [0.99; 0.01]$
 \Rightarrow suppose we get $\underline{l^1}$

Correctness Proof of SAMPLE-LW

- ▶ Think of $\mathcal{X} = \mathbf{Y} \cup \mathbf{E}$ (the \mathbf{Y} are sampled, $\mathbf{E} = \mathbf{e}$ given)
- ▶ Probability P_{LW} that SAMPLE-LW generates a sample $\mathbf{x} = (\mathbf{y}, \mathbf{e})$ is

$$P_{LW}(\mathbf{y}, \mathbf{e}) = \prod_{Y_i \in \mathbf{Y}} P(y_i \mid \text{pa}(Y_i)) \prod_{E_i \in \mathbf{E}} 1.0 = \prod_{Y_i \in \mathbf{Y}} P(y_i \mid \text{pa}(Y_i))$$

- ▶ Weight of sample $\mathbf{x} = (\mathbf{y}, \mathbf{e})$ is

$$w(\mathbf{y}, \mathbf{e}) = \prod_{E_i \in \mathbf{E}} P(e_i \mid \text{pa}(E_i))$$

- ▶ **Weighted probability** of a sample¹ is

$$\begin{aligned} \underline{P_{LW}(\mathbf{y}, \mathbf{e}) \times w(\mathbf{y}, \mathbf{e})} &= \prod_{Y_i \in \mathbf{Y}} P(y_i \mid \text{pa}(Y_i)) \prod_{E_i \in \mathbf{E}} P(e_i \mid \text{pa}(E_i)) \\ &= \prod_{X_i \in \mathcal{X}} P(x_i \mid \text{pa}(X_i)) \\ &= \underline{P(\mathbf{y}, \mathbf{e})} \propto \underline{P(\mathbf{y} \mid \mathbf{e})} \quad \square \end{aligned}$$

¹ A sample influences the final probability estimate for the query through the probability with which it will appear, and through its weight.

Likelihood Weighting: Pros and Cons

Advantages:

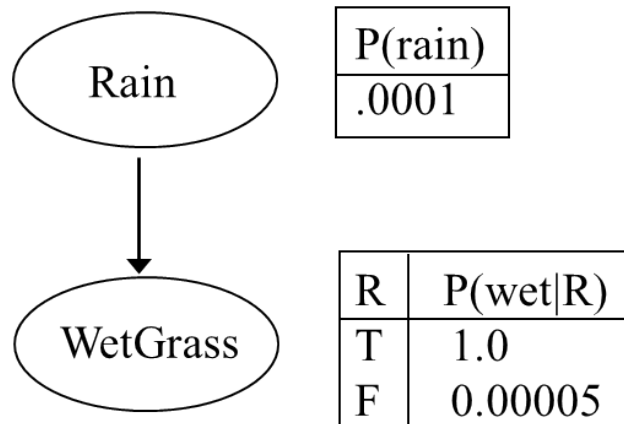
- ▶ Consistent: produces provably correct estimate
- ▶ Less wasteful than Rejection Sampling:
Generates only samples consistent with the evidence.

Disadvantages / Problems:

- ▶ Evidence affects sampling only for nodes *below* the E variables ...
- ▶ ... but **not** for nodes *above* the evidence
(more generally: for any non-descendants of the E)
- ▶ In above example: still sample I ignoring the fact that $S = s^1$
- ▶ Consequence: many samples will be nearly irrelevant to the given evidence;
(correctly) reflected in tiny weights
- ▶ Only few samples typical for the evidence will be available for estimation
- ▶ Particularly bad when most evidence variables are ‘downstream’ in the net
(which is often the case in causally directed networks).

An Extreme Example

Consider this model of the climate in Death Valley:



... and the query

$$P(Rain \mid wet) = ?$$

- ▶ LW will start by sampling *Rain*
- ▶ Might produce thousands of samples $(\neg rain, wet)$
(and correctly attach low weight 0.00005 to each of them)
... before it happens to generate even *one* sample $(rain, wet)$
(which is much more typical given the evidence: weight $w(rain, wet) = 1.0$)
- ▶ In fact, when drawing 5000 random samples, we might not get *even a single* sample $(rain, wet)$!

A More General View: Importance Sampling

General Goal:

- ▶ Draw samples from some distribution P

Problem:

- ▶ May be impossible (or undesirable) to *sample* from P directly
- ▶ But: may be possible to *evaluate* P pointwise (compute $P(x)$ for a given x)

Idea of Importance Sampling:

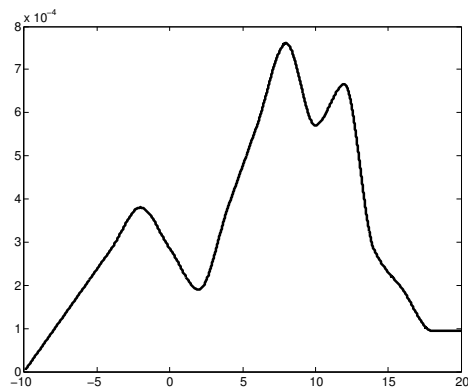
- ▶ Generate samples from a different ('easier') distribution Q instead of P
- ▶ Correct for difference between Q and P by attaching **weights** to samples: high weights where Q under-estimates P , low weights where Q over-estimates P

Terminology:

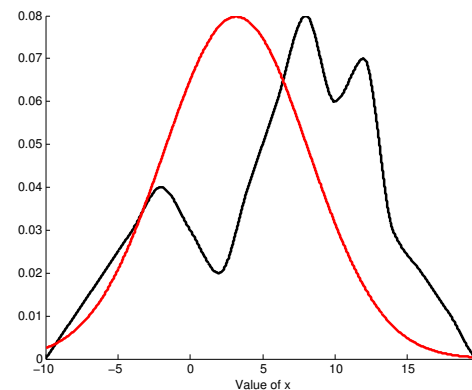
- ▶ P is called **Target Distribution**
- ▶ Q is called **Proposal Distribution**

Importance Sampling

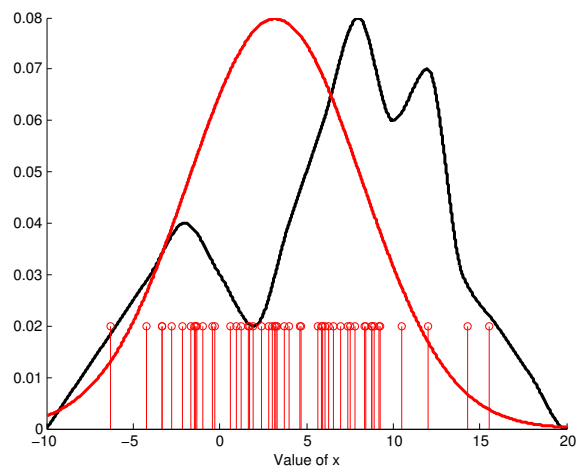
Target distribution P



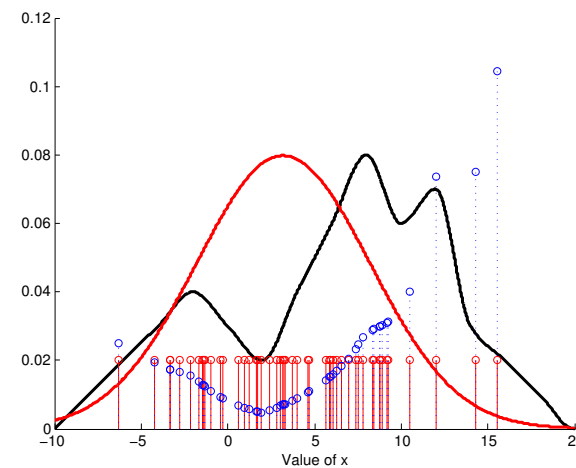
Proposal distribution Q



Points sampled from Q



Weights to correct difference $P : Q$



Importance Sampling

Computing the Weight:

- ▶ $w(x_i)$ of sample x_i must correct for the fact that x_i was drawn from the ‘wrong’ distribution Q
- ▶ Correct weight turns out to be

$$w(x_i) = \frac{P(x_i)}{Q(x_i)}$$

Proof:

- ▶ omitted (but simple)

Intuition:

- ▶ If at some point x the probability/density according to P is n times higher than Q (from which the samples were drawn), then P would generate n times as many samples around x as would Q
- ▶ The Q sample at x must be “upweighted” to count for n samples of P .

Likelihood Weighting as an Instance of Importance Sampling

Reconsider the sampling step SAMPLE-LW in Likelihood Weighting:

Given a query $P(\mathbf{X} \mid e)$, SAMPLE-LW

- ▶ sets all evidence variables E_i to $E_i = e_i$
- ▶ samples the remaining variables $\mathbf{Y} = \mathcal{X} - \mathbf{E}$ in topological order (and weights the resulting sample by $\prod_i P(e_i \mid \text{pa}(E_i))$).

Question:

- ▶ What distribution is SAMPLE-LW actually drawing its samples from?

Definition

Let \mathcal{B} be a network and $\mathbf{E} = e$ fixed evidence ($E_1 = e_1, \dots, E_k = e_k$)

We define the **Mutilated Network** $\mathcal{B}_{\mathbf{E}=e}$ as follows:

- ▶ Variables E_i have *no parents* in $\mathcal{B}_{\mathbf{E}=e}$
- ▶ The CPD of each E_i in $\mathcal{B}_{\mathbf{E}=e}$ gives probability 1.0 to $E_i = e_i$ and probability 0 to all other $e'_i \in \text{Val}(E_i)$
- ▶ The parents and CPDs of all other variables $X \notin \mathbf{E}$ are as in \mathcal{B}

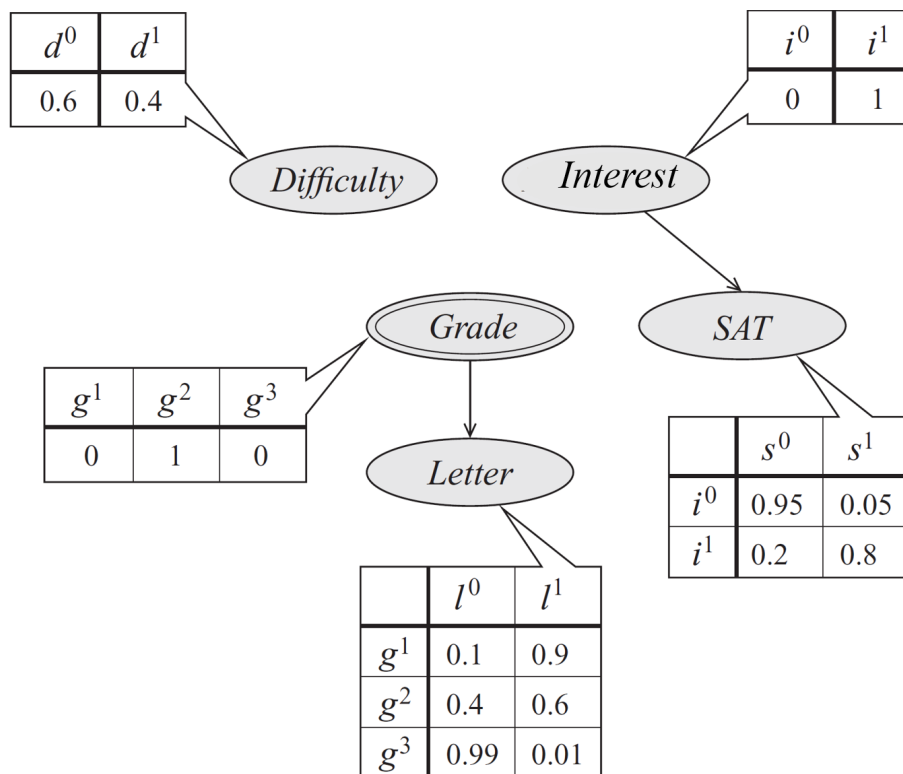
An Example

Query:

$P(D \mid i^1, g^2) = ?$

Corresponding mutilated network

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



Corresponding distribution that we actually sample from:

$$P_{\mathcal{B}_{I=i^1, G=g^2}}(D, S, L) = P(D) \cdot P(S \mid i^1) \cdot P(L \mid g^2) = \prod_{Y_i \in \mathcal{X} - \mathbf{E}} P(y_i \mid \text{pa}(Y_i))$$

Likelihood Weighting as an Instance of Importance Sampling

THEOREM

Let $\langle x, w(x) \rangle$ be a weighted sample generated by SAMPLE-LW.
Then the distribution over x is as defined by the network $\mathcal{B}_{E=e}$ and

$$w(x) = \frac{P_{\mathcal{B}}(x)}{P_{\mathcal{B}_{E=e}}(x)}$$

IN WORDS:

- ▶ Instead of sampling from the full joint distribution $P_{\mathcal{B}}(\mathcal{X})^2$ SAMPLE-LW samples from $P_{\mathcal{B}_{E=e}}(\mathcal{X})$ (which is guaranteed to be consistent with $E = e$)
- ▶ $P_{\mathcal{B}}$ is the **target distribution** P , $P_{\mathcal{B}_{E=e}}$ the **proposal distribution** Q
- ▶ Correct for use of Q to estimate P via weight

$$\boxed{w(x)} = \prod_{E_i \in E} P(e_i \mid \text{pa}(E_i)) = \frac{\prod_{X_i \in \mathcal{X}} P(x_i \mid \text{pa}(X_i))}{\prod_{Y_i \in \mathcal{X} - E} P(y_i \mid \text{pa}(Y_i))} = \frac{P_{\mathcal{B}}(x)}{P_{\mathcal{B}_{E=e}}(x)} = \boxed{\frac{P(x)}{Q(x)}}$$

²which represents the true probabilities over world \mathcal{X} , but which we do not *want to* sample from, because it gives us lots of useless samples with non-matching evidence values

Summary of the Dilemma So Far

The *real* target distribution we would like to sample from (but cannot):

$$P_{\mathcal{B}}(\mathbf{Y} \mid e)$$

Rejection Sampling:

- ▶ samples from full joint distribution $P_{\mathcal{B}}(\mathcal{X})$
(which is correct, but produces many samples that don't match the evidence)
- ▶ has to reduce sample set to \mathcal{D}_e to obtain samples (\mathbf{y}, e) from $P_{\mathcal{B}}(\mathbf{Y} \mid e)$

Likelihood Weighting / Importance Sampling:

- ▶ samples from proposal distribution $P_{\mathcal{B}_{E=e}}(\mathcal{X})$
(which produces only samples with given evidence, but can be very wrong)
- ▶ has to correct for wrong proposal distribution by weights (which can be tiny)

Question:

- ▶ Can't we do better?
- ▶ Can we sample directly from our target $P_{\mathcal{B}}(\mathbf{Y} \mid e)$ (or “almost” ...),
without explicitly calculating it (which would require solving the query)?

Markov Chain Monte Carlo (MCMC) Sampling

General Goal:

- ▶ **Sample directly from $P_{\mathcal{B}}(Y | e)$** (= our *real* target)
- ▶ Then all samples would have the correct e , and would be equally relevant ...

Basic Idea in MCMC:

- ▶ Construct a **sequence of samples** $x^{(0)}, x^{(1)}, x^{(2)}, \dots$ iteratively, with
- ▶ $x^{(0)}$ drawn from some proposal distribution $Q^{(0)}$ and³
- ▶ the distribution $Q^{(i+1)}$ of the next sample $x^{(i+1)}$ depending on the current sample $x^{(i)}$ ⁴
- ▶ Do this such that $Q^{(i+1)}$ is **closer to the target distribution** than $Q^{(i)}$

Consequence:

- ▶ Implicitly, we construct a sequence of proposal distributions $Q^{(0)}, Q^{(1)}, Q^{(2)} \dots$ that **converge** to our target distribution $P_{\mathcal{B}}(Y|e)$
- ▶ Later samples are “almost” drawn from $P_{\mathcal{B}}(Y|e) \Rightarrow$ no weighting necessary!

³ $Q^{(0)}$ might be the mutilated network $\mathcal{B}_{E=e}$

⁴That makes the sequence of distributions a **Markov chain** – see later.

A Simple MCMC Algorithm: Gibbs Sampling

Given:

- ▶ Probabilistic query $P(\mathbf{X} \mid e)$

Goal:

- ▶ Generate a random sample from target distribution $P(\mathbf{Y} \mid e)$

Basic Idea:

- ▶ Start with a sample $\mathbf{x}^{(0)} = (\mathbf{y}^{(0)}, e)$ drawn from mutilated network $\mathcal{B}_{E=e}$ (i.e., use SAMPLE-LW, but ignore weight)
- ▶ Iteratively try to “fix” (improve) the current sample $\mathbf{x}^{(t)}$ by ...
- ... iterating over each non-evidence variable Y_i and **re-sampling** a value for it, **given the current values of all the other variables**
- ▶ Allows information to “flow” across network (from evidence to non-evidence variables).

The Gibbs Sampling Algorithm

Gibbs Sampling Algorithm

Given:

- ▶ Probabilistic query $P(\mathbf{X} \mid \mathbf{e})$
- ▶ Let $\mathcal{X} = \mathbf{Y} \cup \mathbf{E}$

Algorithm:

- ▶ Sample $\mathbf{x}^{(0)} = (\mathbf{y}^{(0)}, \mathbf{e})$ from $\mathcal{B}_{\mathbf{E}=\mathbf{e}}$
- ▶ For $t = 1, \dots, T$ do:
 - $\mathbf{x}^{(t)} \leftarrow \mathbf{x}^{(t-1)}$
 - for each $Y_i \in \mathbf{Y}$:
 - Resample $y_i^{(t)}$ according to $P(Y_i \mid \mathbf{x}_{-y_i}^{(t)})$ % replace value of Y_i in $\mathbf{x}^{(t)}$
- ▶ Return $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}$

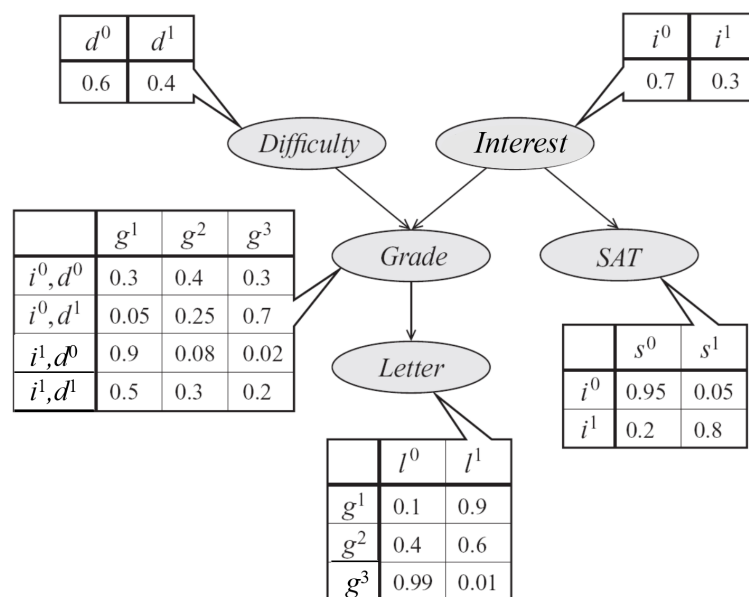
Central Questions:

- ▶ What does the **Resampling Distribution** $P(Y_i \mid \mathbf{x}_{-y_i}^{(t)})$ look like?
- ▶ Can we effectively calculate it (so that we can sample from it)?
- ▶ ... and why does this make sense at all??

Gibbs Sampling: An Example

Query: $P(I \mid s^1, l^0) = ?$

Step 0: Generate first sample: Suppose $x^{(0)} = (d^1, i^0, g^2, \underline{s^1}, \underline{l^0})$



One Gibbs Sampling Iteration:

Loop over $Y = \{G, I, D\}$:

- 1 Resample G acc. to $P(G \mid d^1, i^0, s^1, l^0) \Rightarrow$ suppose we get g^3
- 2 Resample I acc. to $P(I \mid d^1, g^3, s^1, l^0) \Rightarrow$ suppose we get i^1
- 3 Resample D acc. to $P(D \mid i^1, g^3, s^1, l^0) \Rightarrow$ suppose we get d^1

Result: New sample $x^{(1)} = (d^1, i^1, g^3, \underline{s^1}, \underline{l^0})$

Computing the Re-Sampling Distribution

Central Question: How to compute and sample from $P(G \mid d^1, i^0, s^1, l^0)$?

Sampling a single variable, given all others, turns out to be easy:

$$\begin{aligned}
 \boxed{P(G \mid d^1, i^0, s^1, l^0)} &= \frac{P(G, d^1, i^0, s^1, l^0)}{\sum_g P(g, d^1, i^0, s^1, l^0)} \\
 &= \frac{P(d^1) P(i^0) P(G \mid d^1, i^0) P(s^1 \mid i^0) P(l^0 \mid G)}{\sum_g P(d^1) P(i^0) P(g \mid d^1, i^0) P(s^1 \mid i^0) P(l^0 \mid g)} \\
 &= \frac{P(d^1) P(i^0) P(s^1 \mid i^0)}{P(d^1) P(i^0) P(s^1 \mid i^0)} \frac{P(G \mid d^1, i^0) P(l^0 \mid G)}{\sum_g P(g \mid d^1, i^0) P(l^0 \mid g)} \\
 &= \frac{P(G \mid d^1, i^0) P(l^0 \mid G)}{\sum_g P(g \mid d^1, i^0) P(l^0 \mid g)} \\
 &= \boxed{\frac{1}{Z} \times P(G \mid d^1, i^0) P(l^0 \mid G)}
 \end{aligned}$$

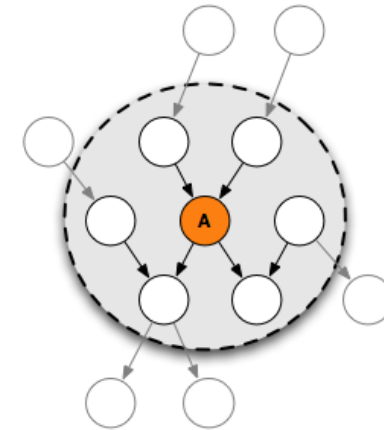
In Words:

- ▶ The re-sampling distribution depends only on G 's parents and G 's child
- ▶ Product of all CPDs in the network that involve G

Computing the Re-Sampling Distribution

Definition

The **Markov Blanket** $Mb(A)$ of a variable A in a Bayesian network is the set of variables consisting of A 's parents, its children, and its children's other parents.



Definition

The **Gibbs Resampling Distribution** for a variable Y_i , given the values of all other variables $\mathcal{X} - Y_i$ in a Bayesian network is defined as:

$$P(Y_i \mid \mathbf{x}_{-y_i}) = P(Y_i \mid mb(Y_i))$$

Easy to show:

- ▶ $P(Y \mid mb(Y))$ is proportional to the probability of Y given its parents, times the probability of each child given its respective parents:

$$P(Y \mid mb(Y)) = \frac{1}{Z} \times P(Y \mid pa(Y)) \times \prod_{C_i \in Children(Y)} P(c_i \mid pa(C_i))$$

Estimating $P(\mathbf{X} | e)$ [3]: Gibbs Sampling

Probabilistic Inference with Gibbs Sampling

Given:

- ▶ A Bayesian network \mathcal{B} over \mathcal{X}
- ▶ A query $P(\mathbf{X} | e)$ with $\mathbf{X}, \mathbf{E} \subseteq \mathcal{X}$

Algorithm:

- 1 **Sample:** Probabilistically sample a set \mathcal{D}_e of M events from \mathcal{B} via Gibbs Sampling^a
- 2 Take \mathcal{D}_e as a direct approximation of the target distribution $P(\mathbf{Y} | e)$
- 3 **Estimate:** For each value combination $\mathbf{X} = \mathbf{x}$, count the number K of samples $\in \mathcal{D}_e$ that match \mathbf{x} . Estimate $P(\mathbf{x} | e)$ as

$$\hat{P}(\mathbf{x} | e) = \frac{K}{|\mathcal{D}_e|} = \frac{K}{M}$$

^aNote: Due to (generally) low dimensionality of the variables' Markov Blankets, we can pre-compute and store all possible re-sampling distributions in tables.

Why Gibbs Sampling Works

Intuitive Explanation:

- ▶ Consider example from before:
Query $P(I \mid \underline{s}^1, \underline{l}^0)$. First sample $\mathbf{x}^{(0)} = (d^1, i^0, g^2, \underline{s}^1, \underline{l}^0)$
- ▶ Create $\mathbf{x}^{(1)}$, first step: Re-sample G according to the resampling distribution

$$P(G \mid d^1, i^0, \underline{s}^1, \underline{l}^0) = \boxed{1/Z \times P(G \mid d^1, i^0) P(\underline{l}^0 \mid G)}$$

⇒ Sampling G now takes into account downstream evidence of G 's child \underline{L}

⇒ Distribution for new instance $\mathbf{x}^{(1)}$ is a bit closer to target distribution

- ▶ Next steps in creating $\mathbf{x}^{(1)}$:
Re-sample I (using G , D and \underline{S}), then re-sample D (using G and I)
- ⇒ Information from evidence nodes propagates through sampling process
- ⇒ Sampling distribution becomes more and more similar to our target distribution $P(\mathbf{Y} \mid \mathbf{e}) \dots$

Formal Explanation:

- ▶ Gibbs Sampling is an **MCMC** algorithm with **stationary distribution**
 $P(\mathbf{Y} \mid \mathbf{e})$

Gibbs Sampling as an MCMC Algorithm

In Words:

- ▶ In Gibbs Sampling, the sequence of samples is produced by a Markov chain where the states are the possible value assignments (samples, atomic events), and the transition model \mathcal{T} is the Gibbs Resampling Distribution
 - ▶ It can be proven⁵ that the **stationary distribution** of this Markov chain is exactly the **posterior distribution** $P(\mathbf{Y} \mid e)$
 - ▶ In other words: As we continue the sampling process for a long time, the samples come from distributions $P^{(t)}$ that converge to our target distribution!
- ▶ Gibbs Sampling is an instance of the general class of MCMC algorithms
 - ▶ If we sample long enough before starting to use the samples, our samples will closely approximate our target distribution $P(\mathbf{Y} \mid e)$.

⁵Proof omitted for time reasons

Markov Chain Monte Carlo (MCMC) Methods

Given:

- ▶ A target distribution P^* we want to estimate (e.g., posterior $P(\mathbf{Y} \mid e)$)

General Idea of MCMC Algorithms:

- ▶ Construct a Markov Chain over the space of possible events (samples) whose sequence of state distributions $P^{(0)}, P^{(1)}, P^{(2)}, \dots$ **converges to the target distribution P^*** (i.e., has P^* as its **stationary distribution**)
- ▶ Draw a sequence of samples $\mathcal{D} = \{x^{(t)}, x^{(t+1)}, x^{(t+2)} \dots\}$ from the chain
- ▶ Estimate target distribution by counting relative frequencies in \mathcal{D}

Central Problem:

- ▶ How to construct a **transition model (re-sampling strategy)** \mathcal{T} that
- ▶ ... guarantees convergence to P^*
- ▶ ... is effectively computable and can be sampled from?

The Metropolis-Hastings Algorithm (briefly)

Metropolis-Hastings is a general algorithm for constructing Markov chains for *any desired Stationary Distribution*

Metropolis-Hastings: Sketch of the Basic Idea

- ▶ General method for constructing a Markov chain \mathcal{T}^P with a particular stationary distribution P^*
- ▶ Unlike Gibbs Sampling, does *not* assume that we can sample from the transition model distribution \mathcal{T}^P directly
- ▶ Rather, uses idea of a *proposal distribution* \mathcal{T}^Q (as in Importance Sampling) as a substitute for \mathcal{T}^P
- ▶ Replace importance weights with *acceptance probabilities* $\mathcal{A}(x \rightarrow x')$: the probability that the next sample proposed by $\mathcal{T}^Q(x^{(t)})$ will be accepted
- ▶ Central part of the algorithm: An analytic method for determining $\mathcal{A}(x \rightarrow x')$ such that the resulting Markov chain has exactly the desired stationary distribution P^* .

Relatively straightforward to formulate M-H for Bayesian Networks
(but not further discussed in this class)

The Metropolis-Hastings Algorithm (briefly)

A Bit of History

The Metropolis-Hastings algorithm originally comes from the field of statistical physics. It was developed by a group of theoretical physicists at the Los Alamos National Laboratory, New Mexico, then the center for classified military research in the U.S. One of the authors of the original publication is Edward Teller, colloquially known as the “father of the hydrogen bomb”. The method was first derived in the 1950’s to approximate the *Boltzmann Distribution* and later generalised by W.K. Hastings in 1970. See:

Metropolis, N.; Rosenbluth, A.W.; Rosenbluth, M.N.; Teller, A.H.; Teller, E. (1953). Equation of State Calculations by Fast Computing Machines. *Journal of Chemical Physics* 21(6): 1087-1092.

Hastings, W.K. (1970). Monte Carlo Sampling Methods Using Markov Chains and Their Applications. *Biometrika* 57(1): 97-109.

Using MCMC: Some Practical Issues

Two (of many other) questions when applying MCMC in practice:

- 1 How long to wait until we start collecting samples from the chain?
(Want to wait until sample distribution $P^{(t)}$ is close to P^*)
- 2 How to sample a set of *uncorrelated* samples from the chain?

Question 1: Burn-in Time

Question 1: Number of steps T until we collect a sample from the chain?

- ▶ “Burn-in Time”

Answer:

- ▶ Depends on how fast (we believe) the chain converges to P^*

Definition

Let \mathcal{T} be a Markov chain; and

T_ϵ the minimal T such that, for any starting distribution $P^{(0)}$,

$$\mathbb{D}(P^{(T)}; P^*) \leq \epsilon$$

(where \mathbb{D} is a measure of distance between distributions).

Then T_ϵ is called the ϵ -Mixing Time of \mathcal{T} .

Mixing Time

- ▶ can be extremely long
- ▶ cannot be calculated analytically
- ▶ is hard to analyse and estimate.

Question 2: How to Collect Samples

Question 2: How to sample a set of *uncorrelated* samples from the chain?

Problem:

- ▶ Consecutive samples will be strongly *correlated*
⇒ cannot be used as *independent samples*
- ▶ Estimation based on these samples is **biased**.

Pragmatic Solution:

- ▶ After collecting a sample $x^{(t)}$, wait for d steps, then collect next sample $x^{(t+d)}$
- ▶ Alleviates problem slightly
- ▶ But: good d again depends on the *mixing time* of the chain ...

Consequence:

- 👉 In practice, simply wait with collecting samples as long as your time budget permits, and choose d as large as you can afford.

Advanced Topics (not treated here)

Reasoning with Collapsed Particles

- ▶ Problem: In high-dimensional spaces, need huge number of samples
- ▶ “Collapsed Particles”: model part of a variable assignment via particles, combined with closed-form representation of distribution over the rest.

Deterministic Search Methods for Approximate Inference

- ▶ Problem: Sampling can waste a lot of effort in low-probability parts of the space
- ▶ Deterministic Search Methods: Explicitly search for high-probability states (samples)

Optimisation Algorithms for Approximate Inference

- ▶ Entirely different approach: Compute approximate solutions incrementally, using search and message-passing algorithms
- ▶ Operate over graphs (clique trees; factor / cluster / region graphs, ...)
- ▶ Large class of methods; lots of literature.

What you should remember of this section

- ▶ Basic idea of approximate reasoning via sampling
- ▶ How to sample values from a discrete distribution
- ▶ How to sample events from an Bayesian Network
- ▶ How to use samples to estimate a (conditional) distribution:
Rejection Sampling, Likelihood Weighting
- ▶ The general idea of Importance Sampling
- ▶ The general idea of MCMC algorithms
- ▶ Gibbs Sampling, and why it is an MCMC method
- ▶ That there is something called Metropolis-Hastings Algorithm
(and what it is for)
- ▶ That there are some non-trivial issues related to actually using this in
practice (mixing time, correlated samples)

Literature

Koller, Daphne and Friedman, Nir (2009).

Probabilistic Graphical Models: Principles and Techniques. Cambridge, MA: MIT Press.

Russell, Stuart J. and Norvig, Peter (2021).

Artificial Intelligence: A Modern Approach (Fourth Edition). Pearson Publishers.