

Likelihood Weighting and Importance Sampling

Sargur Srihari

srihari@cedar.buffalo.edu

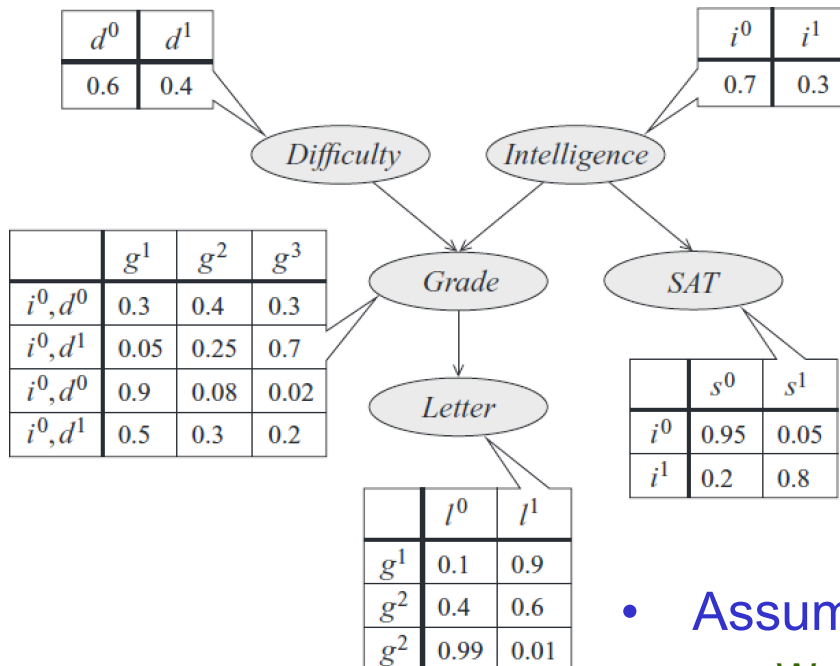
Topics

- Likelihood Weighting Intuition
- Importance Sampling
 - Unnormalized Importance Sampling
 - Normalized Importance Sampling
 - Importance Sampling for Bayesian Networks
 - Mutilated network proposal distribution
 - Computing answers to queries
 - Quality of importance sampling estimator

Inefficiency of Rejection Sampling

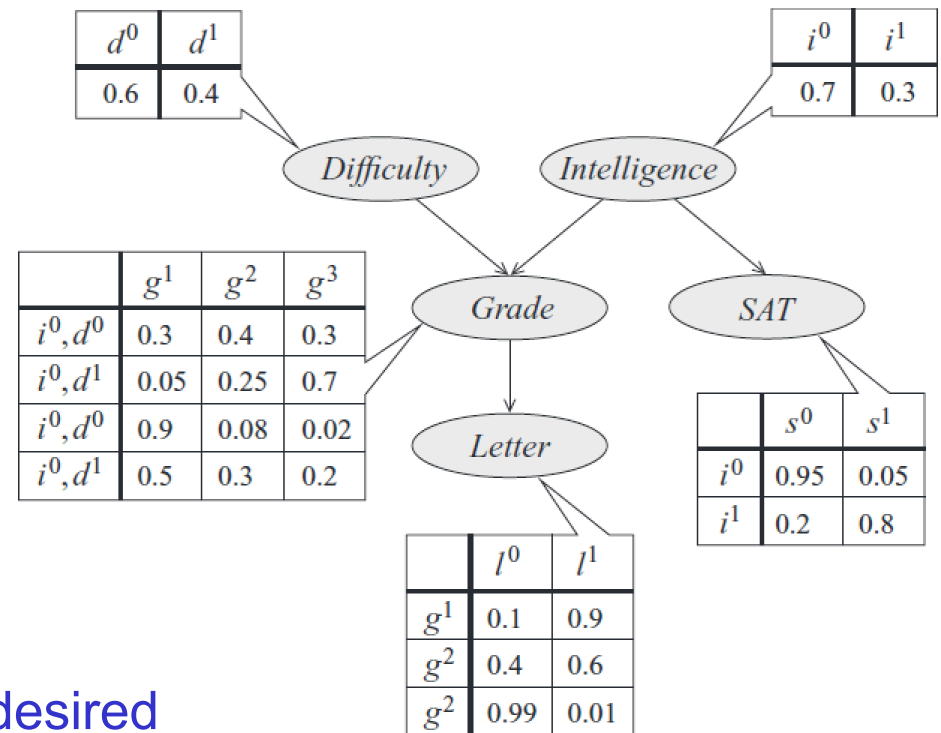
- Rejection sampling process is very wasteful in the way it handles evidence
 - We generate multiple samples that are inconsistent with our evidence
 - They are ultimately rejected without contributing to our estimate
- Here we consider an approach that makes our samples more relevant to our evidence

Inefficiency of Rejection Sampling



- Assume that our evidence is: d^1, s^1
 - We use forward sampling
 - Say it generates: d^0 for D
 - This sample will always be rejected as being incompatible with evidence
- Better approach may be to force samples to match evidence
 - Value takes on only appropriate observed values
 - But it can generate incorrect values as shown next

Forcing Sample Value



- Example of naive forcing process

Evidence is $S=s^1$ (student with high SAT)

1. Using naive process we sample D and I
2. Set $S=s^1$
3. Then sample G and L appropriately

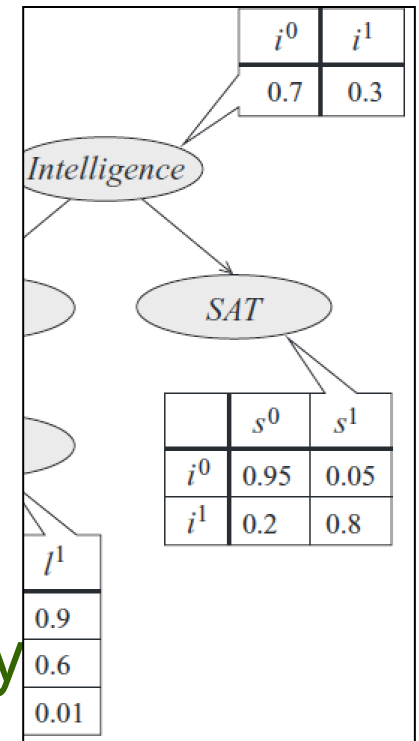
- All of our samples will have $S=s^1$ as desired

- However expected no of samples that have i^1 (an intelligent student) will have probability 30%, same as prior
- Thus it fails to conclude that the posterior of i^1 is higher when we observe s^1

$$P(i^1|s^1)=0.875$$

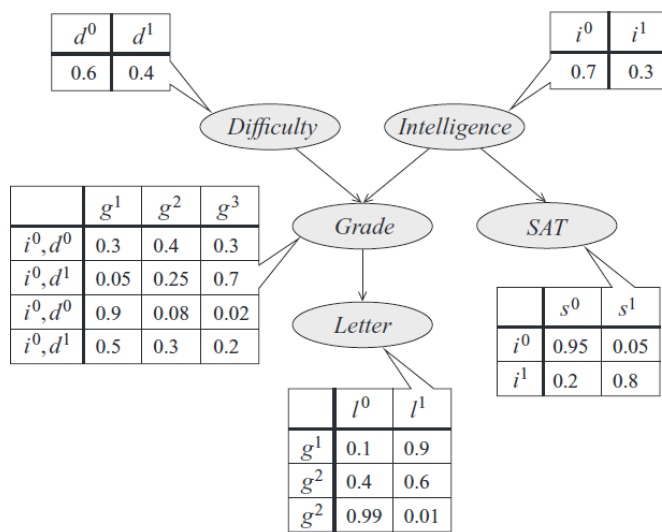
Shortcoming of Forcing

- Forcing $S=s^1$ fails to account for:
 - Node S is more likely to take value $S=s^1$ when parent $I=i^1$ than when $I=i^0$
- Contrast with rejection sampling:
 - With $I=i^1$ we would have generated $S=s^1$ 80% of the time
 - With $I=i^0$ would have generated $S=s^1$ only 5% of the time
- To simulate this long-run behavior with forcing
 - Sample where we have $I=i^1$ and force $S=s^1$ should be 80% of a sample whereas $I=i^0$ and force $S=s^1$ is worth 5% of a sample



Intuition of Weighting

- Weights of samples = likelihood of evidence accumulated during sampling process



- Evidence consists of: ℓ, s^1
- Using forward sampling, assume that we sample $D=d^1, I=i^0$
- Based on evidence, Set $S=s^1$
- Sample $G=g^2$
- Based on evidence, Set $L=\ell$
- Total sample is: $\{D=d^1, I=i^0, G=g^2, S=s^1, L=\ell\}$
- Compensation for evidence ℓ, s^1 which are not sampled
 - Given $I=i^0$, forward sampling would have generated $S=s^1$ with probability 0.05
 - Given $G=g^2$ forward sampling would have generated $L=\ell$ with probability 0.4
- Each event is the result of an independent coin toss
 - probability that both would have occurred is their product
- Weight required for this sample to compensate for the setting for the evidence is $0.05 \times 0.4 = 0.02$

Generalizing this intuition leads to the Likelihood Weighting (LW) algorithm

Algorithm: Likelihood Weighting (LW)

- Algorithm LW-Sample is shown next
- Here weights of samples derived from likelihood of evidence accumulated during sampling process
- This process generates a weighted particle
- It is used M times to generate a set \mathcal{D} of weighted particles

Likelihood-weighted particle generation

- **Procedure** LW-Sample (

\mathcal{B} , // Bayesian Network over χ

$Z=z$ // Event in the network)

- Let X_1, \dots, X_n be a topological ordering of χ

$w \leftarrow 1$

- **for** $i=1, \dots, n$

$u_i \leftarrow \mathbf{x}(\text{Pa}X_i)$ // u_i is assigned parents of X_i among x_1, \dots, x_{i-1}

if $X_i \notin Z$ **then**

- Sample x_i from $P(X_i | u_i)$

else

- $X_i \leftarrow z(X_i)$ // Assignment to X_i in z

- $w \leftarrow w \cdot P(x_i | u_i)$ // Multiply weight by probability of desired value

- **return** $(x_1, \dots, x_n), w$

Process generates
a single weighted
particle given event
 $Z=z$

Likelihood-Weighted (LW) Particles

- Using LW sample to estimate conditional $P(\mathbf{y} | \mathbf{e})$
 - Use it M times to generate a set \mathcal{D} of weighted particles $(\xi[1], w[1]), \dots, (\xi[M], w[M])$
 - We then estimate
$$\hat{P}_D(\mathbf{y} | \mathbf{e}) = \frac{\sum_{m=1}^M w[m] I\{\mathbf{y}[m] = \mathbf{y}\}}{\sum_{m=1}^M w[m]}$$
 - Generalizes expression for unweighted particles with forward sampling
$$\hat{P}_D(\mathbf{y}) = \frac{1}{M} \sum_{m=1}^M I\{\mathbf{y}[m] = \mathbf{y}\}$$
 - Each particle had weight 1; hence terms of numerator unweighted; denominator is sum of all particle weights which is M
 - As in forward sampling same set of samples can be used to estimate $P(\mathbf{y})$

LW: a case of Importance Sampling

- We have not yet provided a formal justification for the correctness of LW
- It turns out that LW is a special case of a very general approach: *Importance Sampling*
- It also provided the basis for analysis
- We first give general description and analysis of importance sampling
- Then reformulate LW as a special case of the framework

Importance Sampling

- Let \mathbf{X} be a set of variables that take on values in some space $Val(\mathbf{X})$
- Importance sampling is a general approach for estimating the expectation of a function $f(\mathbf{x})$ relative to some distribution $P(\mathbf{X})$ typically called the *target distribution*
- We can estimate this expectation by generating samples $\mathbf{x}[1], \dots, \mathbf{x}[M]$ from P

Sampling from a Different Distribution

- Given samples $\mathbf{x}[1], \dots, \mathbf{x}[M]$ we can estimate expectation of f relative to P by

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

- We might prefer to generate samples from a different distribution Q known as the
 - *proposal distribution or sampling distribution*
- Reasons to sample from different distribution:
 - It may be impossible or Computationally very expensive to sample from P
 - Ex: Posterior distribution of a MN or Prior of a MN

Requirements of Proposal Distribution

- We discuss how we might obtain estimates of an expectation relative to P by generating samples from a different distribution Q
- In general proposal distribution can be arbitrary
- We require only that $Q(\mathbf{x}) > 0$ whenever $P(\mathbf{x}) > 0$
 - So that Q does not ignore states that have nonzero probability
 - Support of Q contains the support of P
 - Computational performance depends on the extent to which Q is similar to P

Unnormalized Importance Sampling

- If we generate samples from Q instead of P we cannot simply average the f -value of the samples generated
- We need to adjust our estimator to compensate for the incorrect sampling distribution
- An obvious way of adjusting our estimator is:

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

Proof:

$$\begin{aligned} E_{Q(\mathbf{X})}\left[f(\mathbf{X})\frac{P(\mathbf{X})}{Q(\mathbf{X})}\right] &= \sum_x Q(\mathbf{x})f(\mathbf{x})\frac{P(\mathbf{x})}{Q(\mathbf{x})} \\ &= \sum_x f(\mathbf{x})P(\mathbf{x}) \\ &= E_{P(\mathbf{X})}[f(\mathbf{X})] \end{aligned}$$

Unnormalized Importance Sampling

- Based on the above observation,
 - We generate samples $\mathcal{D}=\{\mathbf{x}[1],\dots,\mathbf{x}[M]\}$ from Q and estimate
$$\hat{E}_D[f] = \frac{1}{M} \sum_{i=1}^M f(\mathbf{x}[m]) \frac{P(\mathbf{x}[m])}{Q(\mathbf{x}[m])}$$
- We call this estimator the unnormalized importance sampling estimator
 - Also called unweighted importance sampling
- The factor $P(\mathbf{x}[m]) / Q(\mathbf{x}[m])$ is a correction weight to the term $f(\mathbf{x}[m])$
- We use $w(\mathbf{x})$ to denote $P(\mathbf{x}) / Q(\mathbf{x})$

Unnormalized Sampling needs P

- Preceding discussion assumed that P is known
- One of the reasons why we must resort to sampling from a different distribution Q is that P is known only up to a normalizing constant Z
 - Specifically we have access to $\tilde{P}(\mathbf{X})$ such that $\tilde{P}(\mathbf{X})$ is not a normalized distribution but $\tilde{P}(\mathbf{X}) = ZP(\mathbf{X})$
 - For example, in a BN \mathcal{B} we might have $P(\mathbf{X})$ is our posterior distribution $P_{\mathcal{B}}(\mathbf{X}|\mathbf{e})$ and $\tilde{P}(\mathbf{X})$ be the unnormalized distribution $P_{\mathcal{B}}(\mathbf{X}, \mathbf{e})$
 - In a MN $P(\mathbf{X})$ might be $P_{\mathcal{H}}(\mathbf{X})$ and $\tilde{P}(\mathbf{X})$ might be product of clique potentials without normalization

Importance Sampling (Normalized)

- Since we do not know P we define

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

- But we cannot use $\hat{E}_D[f] = \frac{1}{M} \sum_{i=1}^M f(\mathbf{x}[m]) \frac{P(\mathbf{x}[m])}{Q(\mathbf{x}[m])}$
- We can use a slightly different estimator based on

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

Its expectation is simply Z

- We can rewrite

$$E_{P(\mathbf{X})}[f(\mathbf{X})] = E_{Q(\mathbf{X})} \left[f(\mathbf{X}) \frac{P(\mathbf{X})}{Q(\mathbf{X})} \right] \quad \text{as:}$$

$$\begin{aligned} E_{P(\mathbf{X})}[f(\mathbf{X})] &= \sum_{\mathbf{x}} P(\mathbf{x}) f(\mathbf{x}) = \sum_{\mathbf{x}} Q(\mathbf{x}) f(\mathbf{x}) \frac{P(\mathbf{x})}{Q(\mathbf{x})} \\ &= \frac{1}{Z} \sum_{\mathbf{x}} Q(\mathbf{x}) f(\mathbf{x}) \frac{\tilde{P}(\mathbf{x})}{Q(\mathbf{x})} = \frac{1}{Z} E_{Q(\mathbf{X})}[f(\mathbf{x}) w(\mathbf{X})] = \frac{E_{Q(\mathbf{X})}[f(\mathbf{x}) w(\mathbf{X})]}{E_{Q(\mathbf{X})}[w(\mathbf{X})]} \end{aligned}$$

- Can use empirical estimator for denominator and num.
- Given M samples $\mathcal{D} = \{\mathbf{x}[1], \dots, \mathbf{x}[M]\}$ from Q

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

This is called the normalized importance sampling estimator

Importance Sampling for BNs

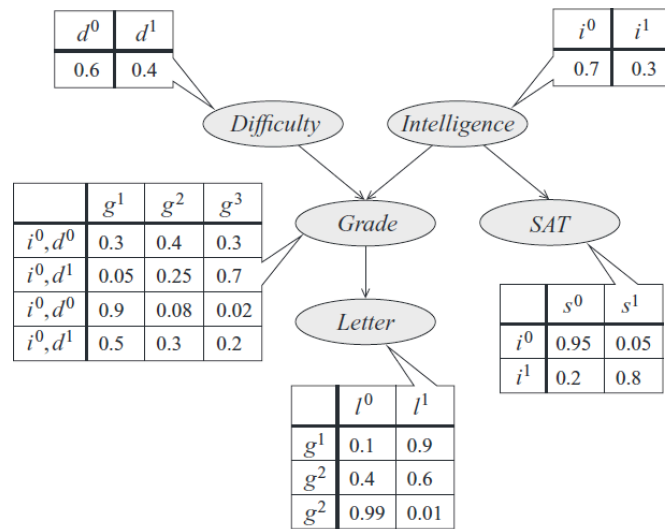
- With the theoretical foundation we can describe application of importance sampling to BNs
- Distribution Q uses network structure/CPDs to focus on part of the joint distribution– the one consistent with a particular event $Z=z$
- Several ways this construction can be applied to BN inference, dealing with various types of probability queries
- Finally discuss several proposal distributions
 - more complex to implement but perform better

Defining the Proposal Distribution

- Assume we are interested in a particular event $Z=z$ either because
 - We wish to estimate its probability, or
 - We have observed it as evidence
- We wish to focus our sampling process on the parts of the joint that are consistent with this event
- We define a process that achieves this goal

Goal of Proposal Distribution

- We are interested in the student's grade $G=g^2$
- We wish to bias our sampling toward parts of the space where this event holds



Easy to take this into account while sampling L

We simply sample from $P(L|g^2)$

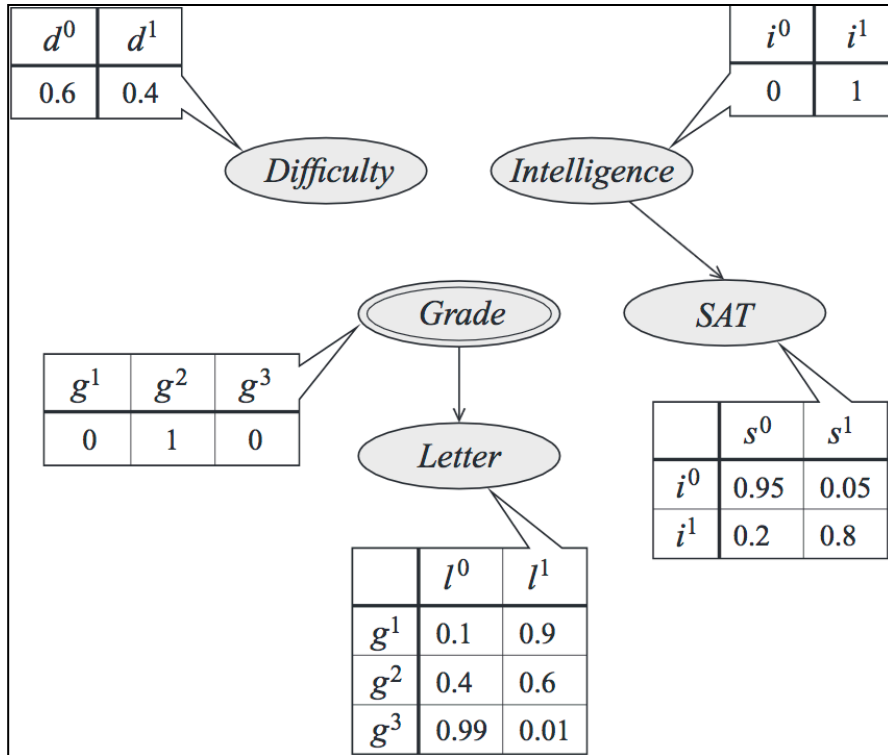
However it is difficult to account for G 's influence on D , L and S without doing inference in the network.

- Goal is to define a simple proposal distribution that allows the efficient generation of particles

Mutilated Bayesian Network, $\mathcal{B}_{Z=z}$

- We define a proposal distribution that sets Z to take pre-specified value in a way that influences the sampling process for its descendants but not the other nodes in the network
- Proposal distribution is described as a BN
 - Each node $Z_i \in Z$ has no parents in $\mathcal{B}_{Z=z}$
 - Parents and CPDs of all other nodes $X \notin Z$ are unchanged

Ex: Mutilated Net for $\mathcal{B}_{I=i1, G=g2}$



- Node G is decoupled from its parents, eliminating its dependence on them
- I has no parents in original network
- Both I and G are deterministic ascribing probability 1 to their observed values

Importance sampling with this proposal distribution is precisely equivalent to the *Likelihood Weighting algorithm* seen earlier

If ξ is a sample generated by the LW algorithm and w is its weight. Then

$$w(\xi) = \frac{P_B(\xi)}{P_{B_{Z=z}}(\xi)}$$

Use of the proposal distribution

- Mutilated net can be used for estimating a variety of BN queries
 1. Unconditional probability of an event $Z=z$
 2. Conditional probability of an event $P(y|e)$ for a specific event y
 3. An entire joint distribution $P(\mathbf{Y}|e)$ for a subset of variables \mathbf{Y}

Computing Unconditional of Event

- Simple problem of computing the unconditional probability of an event $Z=z$
 - We can clearly use forward sampling
 - We can also use unnormalized importance sampling where P is defined by $P_{\mathcal{B}}(X)$ and Q is defined by the mutilated network $\mathcal{B}_{Z=z}$
- Our goal is to estimate expectation of a simple function f which is the indicator function of a query z : $f(\xi) = I\{\xi(Z) = z\}$
 - The estimator for this case is simply

$$P_D(z) = \frac{1}{M} \sum_{m=1}^M I\{\xi[m](Z) = z\} w(\xi[m]) = \frac{1}{M} \sum_{m=1}^M w[m]$$

Computing the Conditional $P(\mathbf{y}|\mathbf{e})$

- Compute it as $P(\mathbf{y}, \mathbf{e}) / P(\mathbf{e})$
 - Called *Ratio Likelihood Weighting (Ratio LW)*
- We use unnormalized importance sampling for both numerator and denominator
 - Estimate conditional probability in two phases
 - Use LW algorithm twice,
 - first M times with $\mathbf{Y}=\mathbf{y}, \mathbf{E}=\mathbf{e}$ to generate \mathcal{D} weighted samples $\{ \xi[1], w[1], \dots, \xi[M], w[M] \}$
 - and then M' times with argument $\mathbf{E}=\mathbf{e}$ to generate \mathcal{D}' samples with weights w'
 - We can then estimate

$$\hat{P}_D(\mathbf{y} | \mathbf{e}) = \frac{\hat{P}_D(\mathbf{y}, \mathbf{e})}{\hat{P}_{D'}(\mathbf{e})} = \frac{1 / M \sum_{m=1}^M w[m]}{1 / M' \sum_{m=1}^{M'} w'[m]}$$

Computing the Conditional $P(\mathbf{Y}|e)$

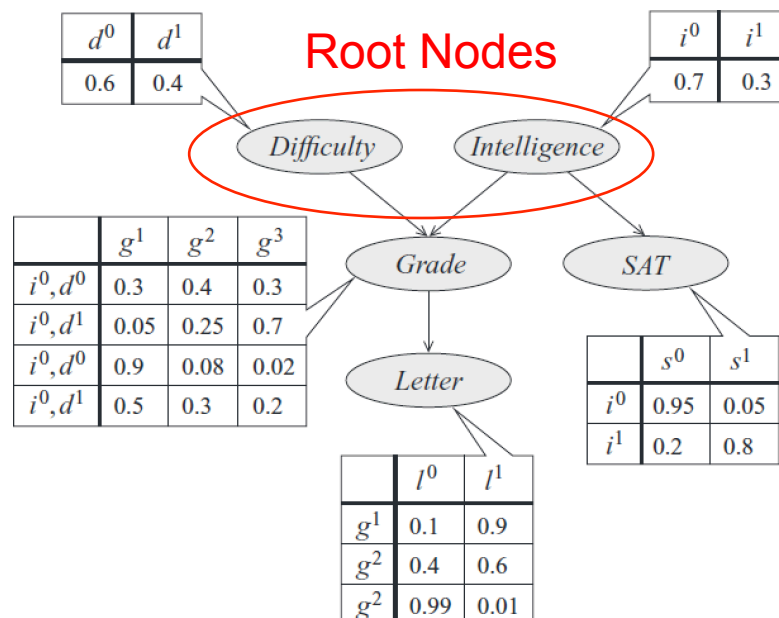
- For a subset of variables \mathbf{Y}
- Running Ratio LW for each $\mathbf{y} \in \text{Val}(\mathbf{Y})$ is computationally impractical
- An alternative approach is to use is to use normalized likelihood weighting

Quality of Importance Sampling

- Depends on how close the proposal distribution Q is to the target distribution P .
- Extreme cases:
 - All evidence at root nodes
 - All evidence is at leaf nodes
- Discussed next

All evidence at root nodes

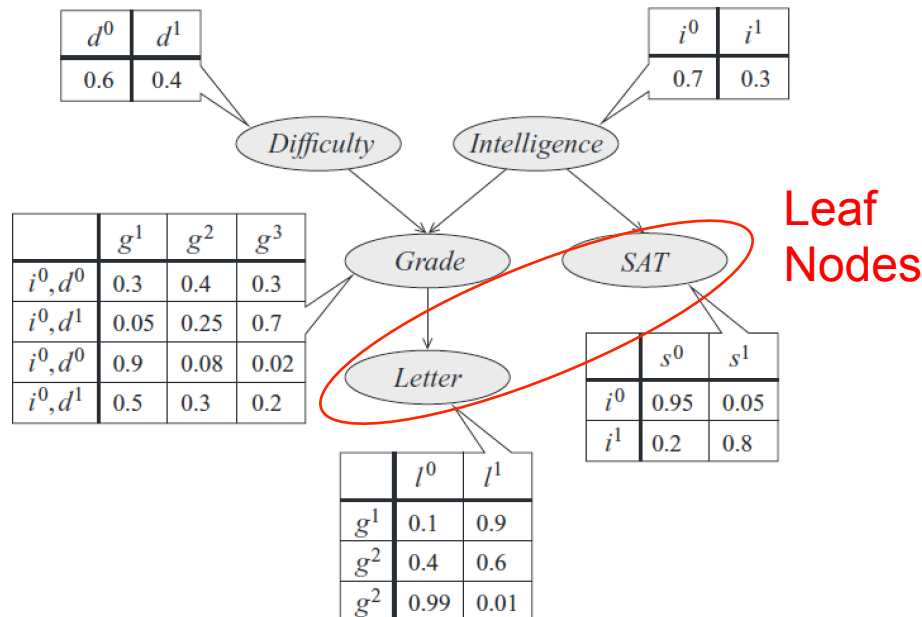
- Proposal distribution is precisely the posterior
- No evidence encountered along the way
- All samples will have the same weight $P(e)$



Q : Posterior: $P(G, S, L | D = d^1, I = i^1)$

All evidence is at leaf nodes

- Proposal distribution $Q(\chi)$ is the Prior distribution $P_{\mathcal{B}}(\chi)$, leaving the correction purely to the weights
- Will work well only if prior is similar to posterior
- Otherwise most samples will be irrelevant



Q : Prior: $P(G, S, L, D, I)$