

Graphical Models

Lecture 15:

Approximate Inference by Sampling

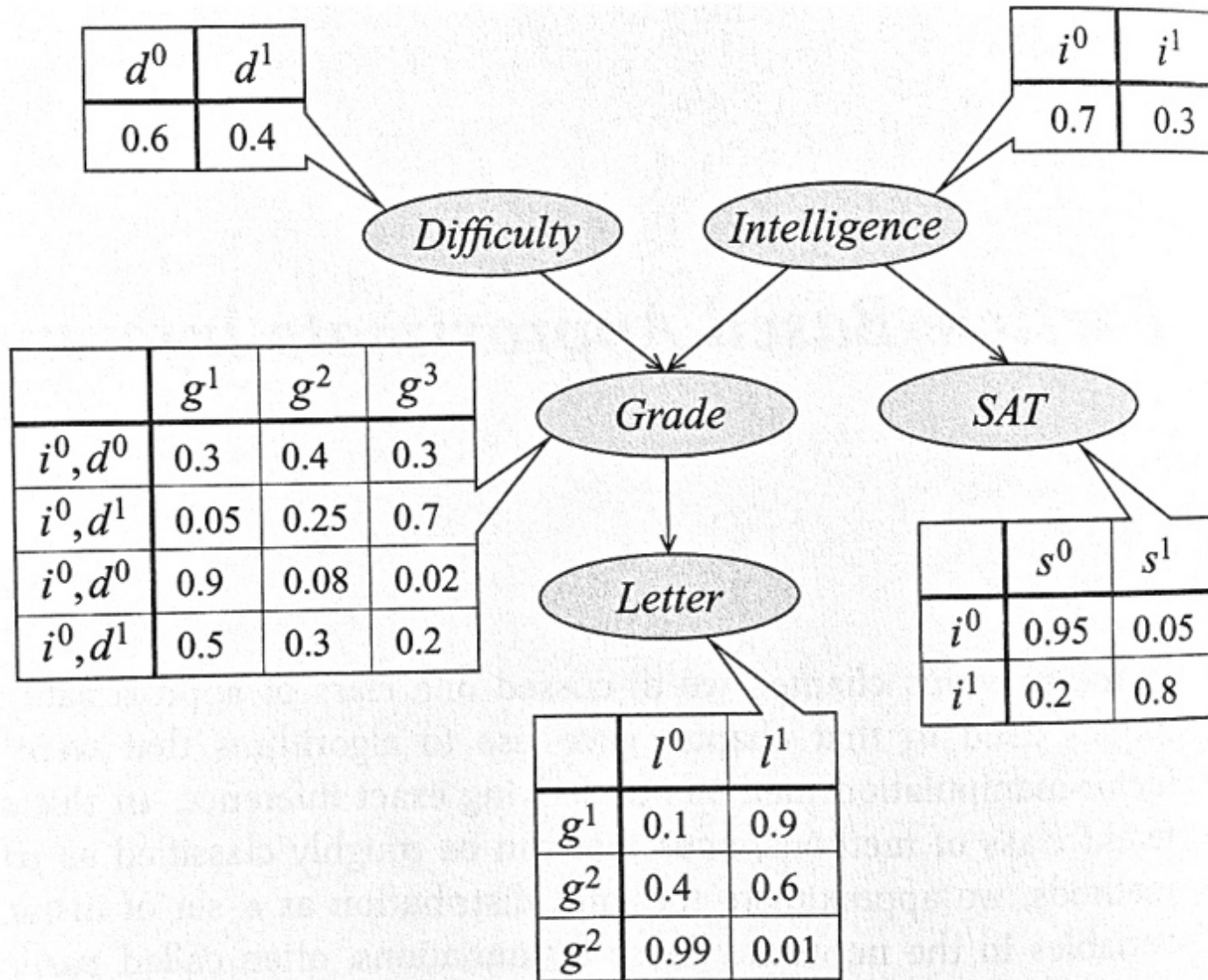
Andrew McCallum
mccallum@cs.umass.edu

Thanks to Noah Smith and Carlos Guestrin for some slide materials.

General Idea

- Set of random variables \mathbf{X} , distribution P
- We want to estimate $E[f(\mathbf{X})]$
 - Often this is the value of some X_i (i.e. a marginal)
 - Could do this by enumerating *all* \mathbf{X} , as you did in HW#1. Instead...
$$\sum_{\mathbf{x}} P(\mathbf{X} = \mathbf{x}) f(\mathbf{x})$$
- Generate M “particles” (samples from P)
 - Calculate f on each one
$$\hat{\mathbb{E}}[f] = \frac{1}{M} \sum_{m=1}^M f(\mathbf{x}^{(m)})$$
 - Aggregate
$$\hat{P}(\mathbf{y}) = \frac{1}{M} \sum_{m=1}^M \mathbb{I}\{\mathbf{y}^{(m)} = \mathbf{y}\}$$
- The methods are *randomized*, unlike last week.

Example



Forward Sampling from a Bayesian Network

- Sample all random variables \mathbf{X} by traversing the Bayesian network in topological order
- Sample each X_i given the (already sampled) values of its parents

- Repeat M times.

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

- Aside: how to sample from a multinomial distribution?

How Many Samples?

- More samples, intuitively, give a better estimate.
- The event for each $\mathbf{y}^{(m)}$ taking the value \mathbf{y} can be understood as a Bernoulli trial.
 - M independent trials
- Hoeffding bound: $P\left(\hat{P}(\mathbf{y}) \notin [P(\mathbf{y}) - \epsilon, P(\mathbf{y} + \epsilon)]\right) \leq 2e^{-2M\epsilon^2}$
- To get absolute error bounded by ϵ with probability $> 1 - \delta$, need $M \geq \frac{\ln \frac{2}{\delta}}{2\epsilon^2}$

But absolute error not good for small p .

How Many Samples?

- We can use the Chernoff bound to determine how many samples are required to get within *relative* error ϵ .
- This bound depends on $P(\mathbf{y})$:
$$M \geq 3 \frac{\ln \frac{2}{\delta}}{P(\mathbf{y})\epsilon^2}$$

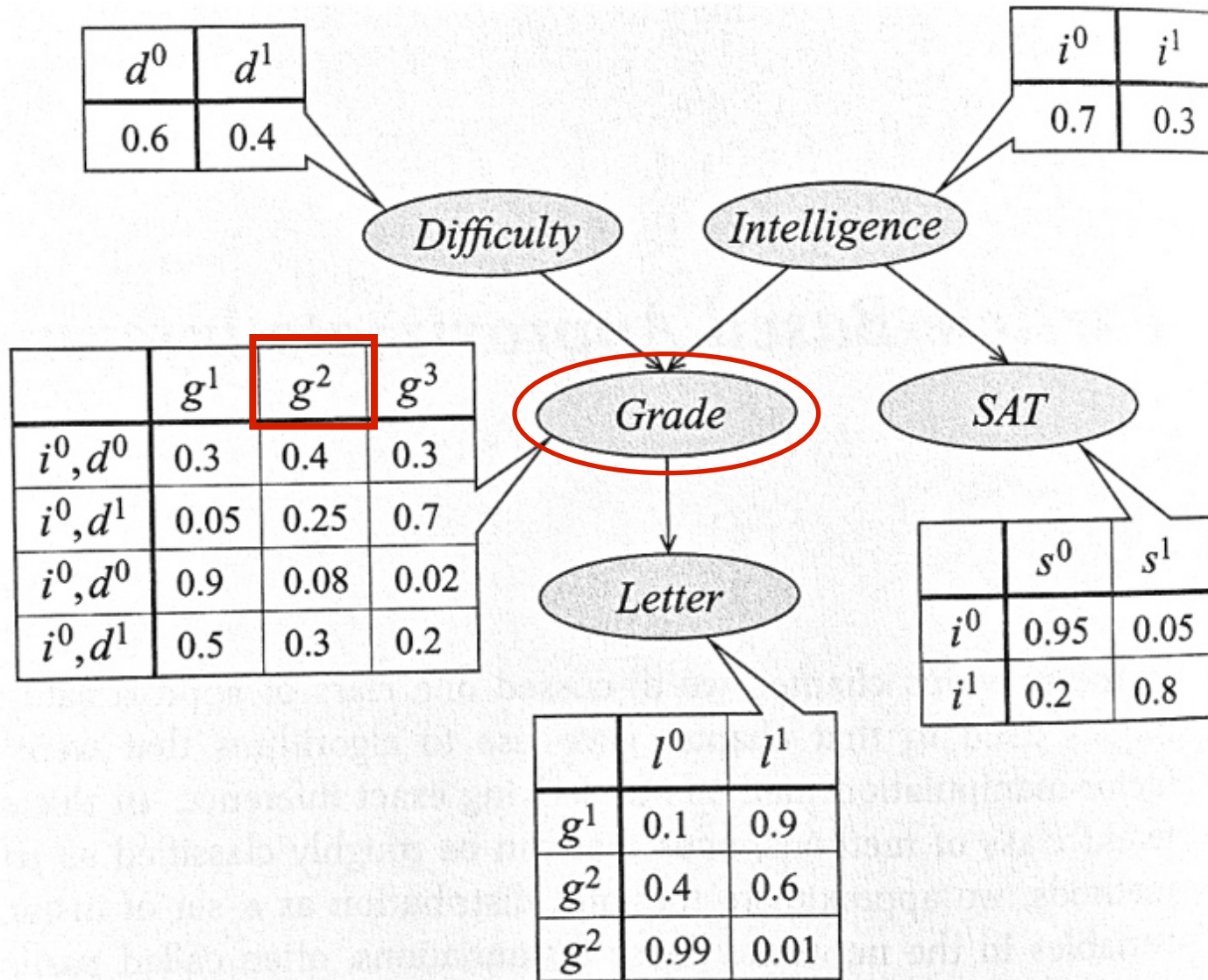
(If we knew $P(\mathbf{y})$ we wouldn't need to estimate it!

We can't determine how many samples we need for a good estimate.)

Dealing with Evidence

- How to use sampling to get $P(\mathbf{Y} | \mathbf{E}=\mathbf{e})$?

Example



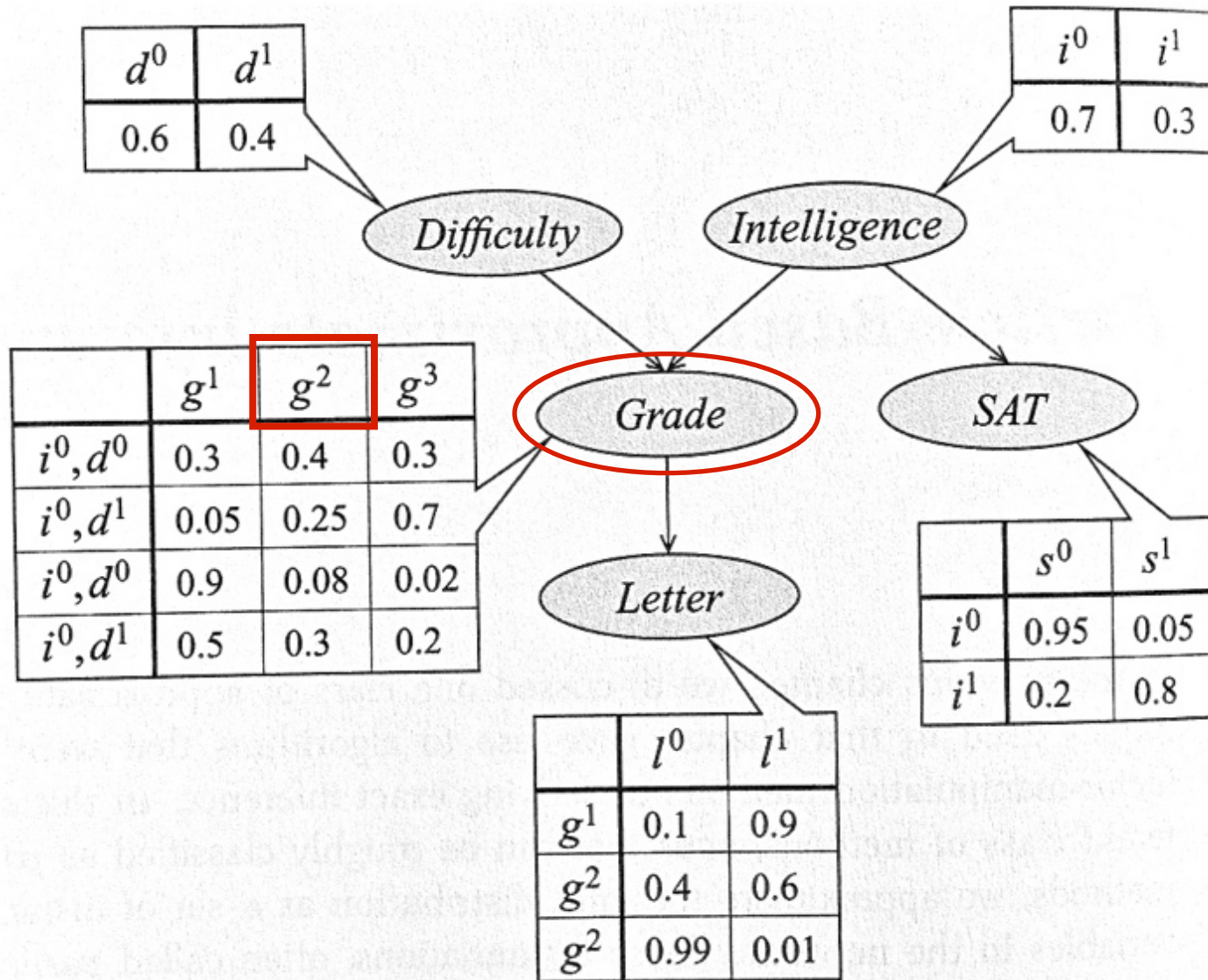
Given Evidence $\mathbf{E} = \mathbf{e}$

- Rejection sampling: draw samples.
 - If inconsistent with evidence, discard.
- Number of required samples is now $M / P(\mathbf{E}=\mathbf{e})$
- Alternatively: estimate $P(\mathbf{Y} = \mathbf{y}, \mathbf{E} = \mathbf{e})$ and $P(\mathbf{E} = \mathbf{e})$ and take their ratio.
 - Low relative error on both will transfer (but that's hard); low absolute error will not.

Likelihood Weighting

- Consider forward sampling, but when we hit an evidence variable, we pick the evidence.
 - Clearly wrong!

Example



Likelihood Weighting

- Basic idea: weight each particle by the probability of the evidence variables taking the values we insist on.
- Estimate:

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

Likelihood Weighting

Input: Bayesian network, evidence \mathbf{e}

- $w = 1$
- Consider each X_i in topological order:
 - If X_i is an evidence variable:
 - $x_i \leftarrow e_i$
 - $w \leftarrow w \cdot P(X_i \mid \text{parents}(X_i))$
 - Else:
 - Sample x_i given $\text{parents}(X_i)$
- Return \mathbf{x} and w

Note

- Forward sampling is a special case of likelihood weighting.
- Next: **importance** sampling generalizes likelihood weighting.

Sampling from a Proposal Distribution

- Let Q be an alternative distribution such that $Q > 0$ whenever $P > 0$. (Idea: easier to sample from Q than from P .)

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

Show this on board.
p. 495

- So given particles, we can weight them by $P(\mathbf{x}^{(m)}) / Q(\mathbf{x}^{(m)})$.

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

Unnormalized Importance Sampling

- (Name is confusing.)
- Expected value of this estimator for f is the true expected value under P .
- Variance of the estimate can be diminished with:
 - Getting more samples
 - Making Q closer to (proportional to): $P \cdot |f|$

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

Normalized Importance Sampling

- We may not be able to calculate $P(\mathbf{x}^{(m)})$.
 - Posterior distribution in a Bayesian network (i.e., we are given evidence)
 - Partition function in a Markov network

- Idea:

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

Yes, *normalized* importance sampling uses the *unnormalized* probability. (Funny naming)

Normalized Importance Sampling

- Estimate the numerator and the denominator.

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

$$w(\mathbf{x}^{(m)}) = \frac{U(\mathbf{x}^{(m)})}{Q(\mathbf{x}^{(m)})}$$

Normalized Importance Sampling

- The weight is itself a random variable.

$$\begin{aligned}\mathbb{E}_Q[w] &= \sum_{\mathbf{x}} Q(\mathbf{x}) \frac{U(\mathbf{x})}{Q(\mathbf{x})} \\ &= \sum_{\mathbf{x}} U(\mathbf{x}) \\ &= Z\end{aligned}$$

Normalized Importance Sampling

- Not unbiased; bias decreases as $1/M$.
- Variance is typically lower than the unnormalized estimator
 - But not always; can construct cases where it's worse.

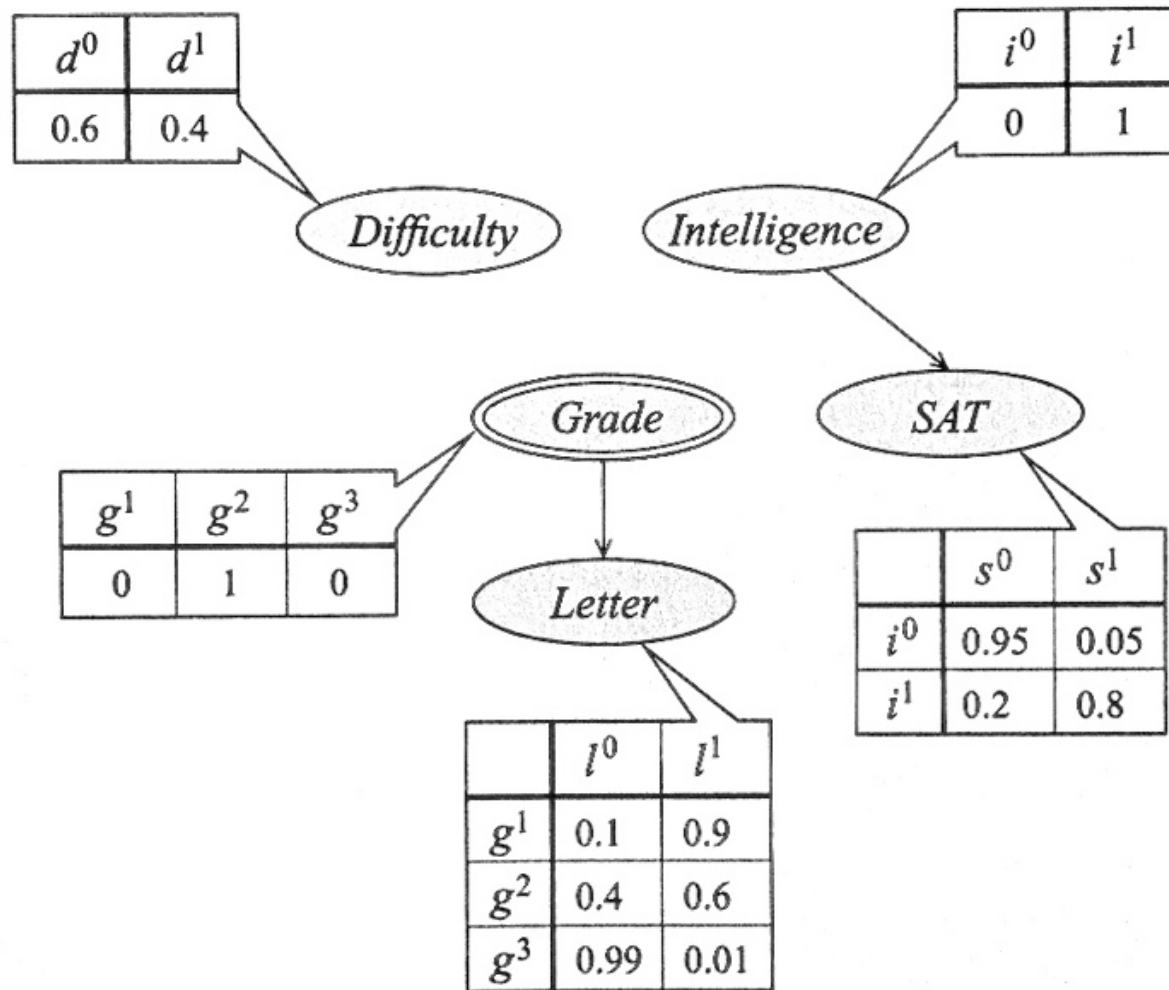
Consider $M=1$,
we get $E_Q[f(\mathbf{X})]$
p 497

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

- We can use the above to measure the effective sample size and tell us whether to keep generating samples.

Next: How to use importance sampling in graphic models?

Mutilated Network





Importance Sampling for Bayesian Networks

- Use Importance Sampling for $P(\mathbf{Y} \mid \mathbf{E}=\mathbf{e})\dots$
- Let the proposal distribution Q be the original network, except:
 - Delete parent connections on evidence nodes.
 - Equivalent to Likelihood Weighting!
- Q called “mutilated network”

Importance Sampling for Bayesian Networks

- Ratio likelihood weighting: calculate the conditional probability of a *specific* event \mathbf{y} : $P(\mathbf{Y} = \mathbf{y} \mid \mathbf{E} = \mathbf{e})$ using two runs of unnormalized importance sampling.
- Normalized likelihood weighting: calculate the conditional probability for a *full distribution* on \mathbf{Y} : $P(\mathbf{Y} \mid \mathbf{E} = \mathbf{e})$ using normalized importance sampling and the LW estimate.


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


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

Quality of Importance Sampling Estimators

- If evidence is at roots, it's perfect.
- If evidence is closer to the leaves, Q is basically the prior over unobserved nodes (ignores the evidence).
- Ratio LW: often lower variance (simply sets the values of \mathbf{Y}); easier to analyze.
- Normalized LW: lets us reuse computation to estimate the whole distribution over \mathbf{Y} .

Importance Sampling

- Very general framework; as long as weights are correct, the process is correct.
 - (In the sense that large enough M gives you a good sample.)
 - Many tricks (e.g., backward importance sampling) for getting a better proposal distribution.

Disadvantages of traditional *Forward Sampling* methods

- Evidence only effects sampling of descendants.
- If evidence is in the leaves of the network, just sampling from the prior.
Could be far from posterior!

Markov Chain Monte Carlo Methods

- Intuition:
“Fix” an initial sample by resampling some of its variables...
- Useful in both directed and undirected models (unlike Likelihood Weighting and other forward sampling methods)

Markov Chain Monte Carlo Methods

- Generate a *sequence* of samples.
- In the beginning, might not be sampling from P .
- Over time, the sampling distribution gets closer and closer to P .

Gibbs Sampling

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

- Sample $\mathbf{x}^{(0)}$ from $P^{(0)}$
- For $t = 1$ to T :
 - $\mathbf{x}^{(t)} \leftarrow \mathbf{x}^{(t-1)}$
 - For each X_i :
 - Sample $x_i^{(t)}$ from $P_{\Phi}(X_i \mid \mathbf{x}_{-i}^{(t)})$
- Return $\mathbf{x}^{(0)} \dots \mathbf{x}^{(T)}$

Notes

- On one round, we sample one random variable given fixed values of all the rest
 - This only depends on some factors
 - Only need to look at some values of those factors
- Appealing: unlike forward sampling, each variable gets to influence values of the others, including accounting for downstream evidence.
- To add evidence: reduce factors at beginning.

Markov Chains

- Graph of states
 - (Do not confuse the states with the graphical model!)
 - One state per element of $\text{Val}(\mathbf{X})$
- Our sampler takes a *random walk* through the states.
- We define the transition probabilities so as to achieve some desirable properties.

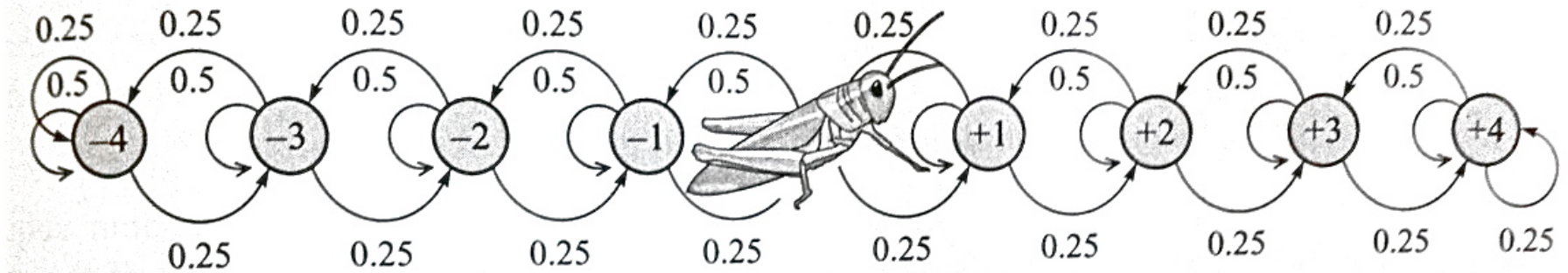
Markov Chains

- Chain dynamics (linear equations):

$$P^{(t+1)}(S^{(t+1)} = s') = \sum_s P^{(t)}(S^{(t)} = s) t(s \rightarrow s')$$

- Long-term behavior: we want the above to converge to $P(\mathbf{X})$.

Drunken Grasshopper



Markov Chains

- **Stationary** distribution:
$$\pi(s') = \sum_s \pi(s)t(s \rightarrow s')$$
- Want a Markov chain that has a *unique* stationary distribution.
 - Problem cases: periodic chains, reducible chains
 - Sufficient condition: **regularity**.
There exists some k such that the probability of getting from any s to any s' in exactly k steps is positive.
 - To get regularity: connectivity, self-loops.

Different Types of Steps

- Our state space will have a factorized structure; we can't move from any state to any other.
- General idea: different transition models can be combined.
 - Select one uniformly at random.
 - Cycle through.
- Though each may not be ergodic, the combination may be.

Two-Dimensional Drunken Grasshopper

Gibbs Sampling and Markov Chains

- Each step of the Gibbs sampler visits a new state in our Markov chain.
- Over time, the sampler converges to the stationary distribution.
- Factors let us calculate each step efficiently.
- By cycling through all random variables, we combine different transition models.

Block Sampling

- Sometimes we can efficiently sample many variables at once.
 - Plate models
- We can think of this as a “big step” in the Markov chain.

Convergence

- It can take a very long time for an MCMC sampler to *mix* (converge to the stationary distribution).

Give an example of a state-transition diagram or graphical model that would not mix well.

- Probably longer than you are willing to wait, for realistic problems.
- We often ignore the earlier elements of the sequence (wait for “burn-in”).
 - Very heuristic.

Generalization: Metropolis-Hastings

- In some situations, even sampling from the posterior for one random variable given all the rest is hard.
- Idea: proposal distribution.
 - Sample X_i from Q , then decide to accept the sample or reject (and stay put).

Metropolis-Hastings

- Still a Markov chain.
- Trade error in the proposal distribution for slower convergence.
- We can show that with the right *accept/reject probabilities*, this Markov chain has the right stationary distribution!

$$A(\mathbf{x}_{-i}, x_i \rightarrow \mathbf{x}_i, x'_i) = \min \left(1, \frac{P_{\Phi}(x'_i, \mathbf{x}_{-i})}{P_{\Phi}(x_i, \mathbf{x}_{-i})} \cdot \frac{t_Q(\mathbf{x}_{-i}, x'_i \rightarrow \mathbf{x}_i, x_i)}{t_Q(\mathbf{x}_{-i}, x_i \rightarrow \mathbf{x}_i, x'_i)} \right)$$

- Gibbs is a special case.

Z's cancel!

MCMC in Practice

- The big problem is mixing time.
- Peakier (more deterministic) distributions make it less probable that the sampler will mix fast.
 - Annealing, or tempering, is a technique that can help.
- Consecutive samples are correlated. Should we wait between drawing particles from the sequence?
 - Probably not, unless f is very expensive to compute. Throwing away particles won't reduce variance.
- Long range moves will speed up mixing time.