

Lecture 17: March 23, 2015
cs 573: Probabilistic Reasoning
Professor Nevatia
Spring 2015

Review

- HW #5 due Mar 30;
 - Conversion from conditional to canonical form posted on class page
- Exam 1, please return with or without comments
- Previous Lecture
 - Variational approximation examples
 - Intro to sampling or particle based approximation
- Today's objective
 - Importance sampling
 - Markov Chain Monte Carlo (MCMC) methods

Likelihood Weighting

- We can just set evidence variables to the observed values to avoid rejecting samples
- However, this biases the samples. Suppose that evidence is $S=s^1$;
- Apply normal forward sampling, we would start from prior distributions of D and I , so $I=i^1$ in only 30% of the cases; this is clearly not consistent with $S=s^1$.
- In rejection sampling, if we started with $I=i^1$ then $S=s^1$ in 80% of the samples (from CPD); if we started with $I=i^0$ then $S=s^1$ in only 5% samples
- Weight samples by their “likelihood”: (i^1, s^1) by .8, (i^0, s^1) by .05
 - This weighting compensation is intuitive, formal justification comes later
 - Generalizes to case where more than one variable is evidence variable (product of probabilities for generating evidence variables given their parents; specifics in Algorithm 12.2, next slide).

$$\hat{P}_{\mathcal{D}}(\mathbf{y} \mid \mathbf{e}) = \frac{\sum_{m=1}^M w[m] \mathbf{I}\{\mathbf{y}[m] = \mathbf{y}\}}{\sum_{m=1}^M w[m]}$$

LW Algorithm

Algorithm 12.2 Likelihood-weighted particle generation

Procedure LW-Sample (
 \mathcal{B} , // Bayesian network over \mathcal{X}
 $Z = z$ // Event in the network
)

- 1 Let X_1, \dots, X_n be a topological ordering of \mathcal{X}
- 2 $w \leftarrow 1$
- 3 **for** $i = 1, \dots, n$
- 4 $u_i \leftarrow x \langle \text{Pa}_{X_i} \rangle$ // Assignment to Pa_{X_i} in x_1, \dots, x_{i-1}
- 5 **if** $X_i \notin Z$ **then**
- 6 Sample x_i from $P(X_i \mid u_i)$
- 7 **else**
- 8 $x_i \leftarrow z \langle X_i \rangle$ // Assignment to X_i in z
- 9 $w \leftarrow w \cdot P(x_i \mid u_i)$ // Multiply weight by probability of desired value
- 10 **return** $(x_1, \dots, x_n), w$

Importance Sampling

- Generalization of LW. Also proves that LW is correct
- It may be hard to generate samples from $P(\mathbf{X})$, called the *target distribution*; instead, sample from another, simpler but related distribution, $Q(\mathbf{X})$, called a *proposal or sampling distribution*.
 - Average the f value (to computed expected value of $f(x)$)
 - This will not give correct expectation of f for P ; need to compensate as shown on next slide.
- Require that $Q(x) > 0$ whenever $P(x) > 0$; otherwise, Q can be arbitrary .
 - *e.g.* $Q(x)$ can be a uniform distribution
 - However, computational efficiency will be better if Q is similar to P

Unnormalized Importance Sampling

- $$\begin{aligned} E_{P(\mathbf{X})} [f(\mathbf{X})] &= \sum_{\mathbf{x}} f(\mathbf{x}) \cdot P(\mathbf{x}) \\ &= \sum_{\mathbf{x}} Q(\mathbf{x}) \cdot f(\mathbf{x}) \cdot (P(\mathbf{x}) / Q(\mathbf{x})) \\ &= E_{Q(\mathbf{X})} [f(\mathbf{X}) P(\mathbf{X}) / Q(\mathbf{X})] \end{aligned}$$

Note: we can compute $P(\mathbf{x})$ for a given \mathbf{x} as the distribution is given

- Given set of samples $D = \{\xi[1], \dots, \xi[M]\}$ from Q

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

Note: last term is like a weight; still this method is called unnormalized importance sampling

- Can be shown that this estimator is not biased: mean is always the desired value
- Variance of the estimator
 - Formal derivation in the book but depends on the variance of $P(\mathbf{x})/Q(\mathbf{x})$, so it is useful to try to find Q that is similar to P .

Normalized Importance Sampling

- Previous method requires computing $P(\xi)$ (normalized) for each sample.
- Sometimes, this is difficult, e.g. in a MN, we can easily get an unnormalized probability but computing partition function may be expensive
- Define a weight relative to the unnormalized distribution

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

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

Normalized Importance Sampling

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

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

Similar to likelihood weighting but samples from Q , not necessarily same as P .

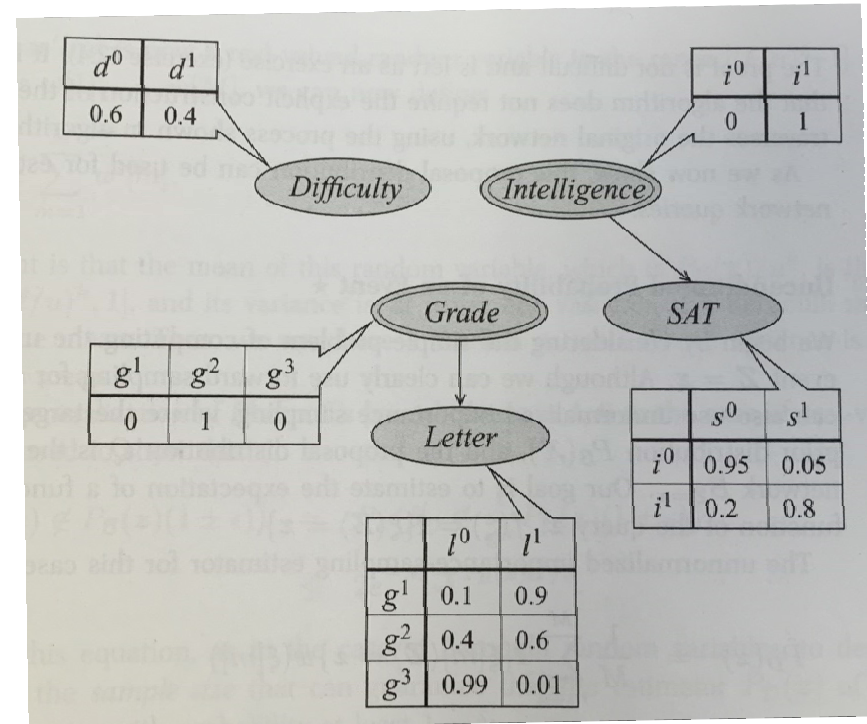
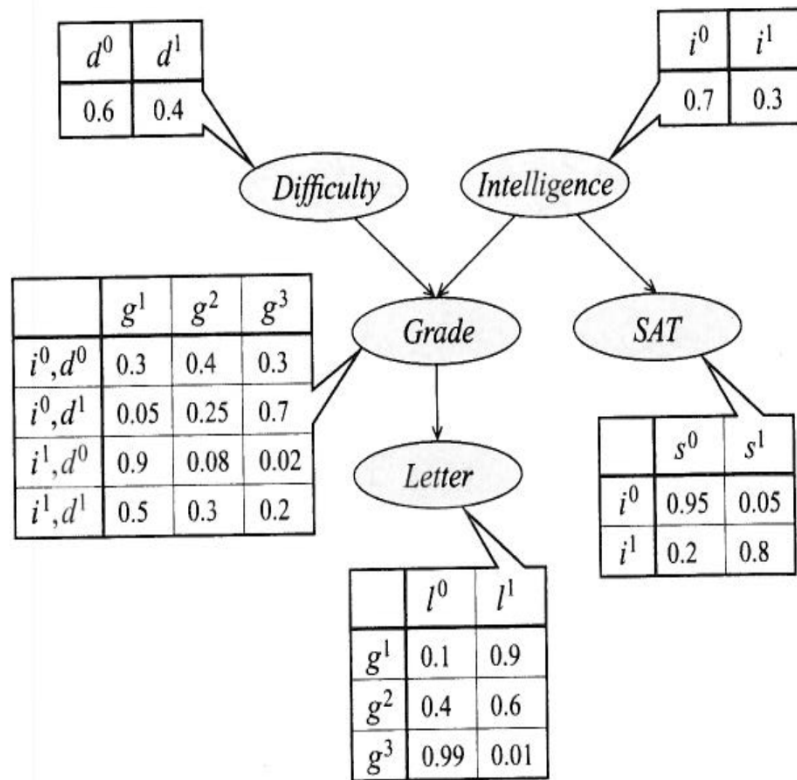
Properties of Normalized IS

- Estimator has a bias for small values of M , estimates are closer to those of Q than of P
 - Bias goes down as $1/M$
- Variance of estimator
 - Analyzed in the book
 - Variance is “typically” lower than for unnormalized sampling though this is not guaranteed
- Normalized sampling may be preferred due to lower variance even though it may have a bias.

Importance Sampling for BN

- For BN, there is a simple construct available for generating a good proposal distribution Q
- Suppose we are interest in the event where $\mathbf{Z} = \mathbf{z}$; e.g. $G = g^2$
 - This could represent evidence or the specific part of distribution we want to estimate
- In example, influence of Z on its descendants (starting from children) is easy to model: e.g. sample from $P(L|g^2)$.
 - Modeling influence on non-descendants is difficult so one choice is just to ignore this influence
- Modify the network so nodes containing Z_i have no parents and CPD of Z_i is 1 for $Z_i = z_i$, 0 otherwise.
- This defines a *mutilated* network (see next slide). Use this network for the proposal distribution function, Q
 - Property that $Q(\mathbf{x}) > 0$ when $P(\mathbf{x}) > 0$ is satisfied

Mutilated Network



Importance Sampling for BN

- Don't actually need to break the network, can incorporate its effect in forward sampling; algorithm 12.2
 - Note: weight modified only when Q and P give different samples
 - Proposition 12.2 shows that the intuitive LW algorithm is equivalent to normalized importance sampling for BNs.

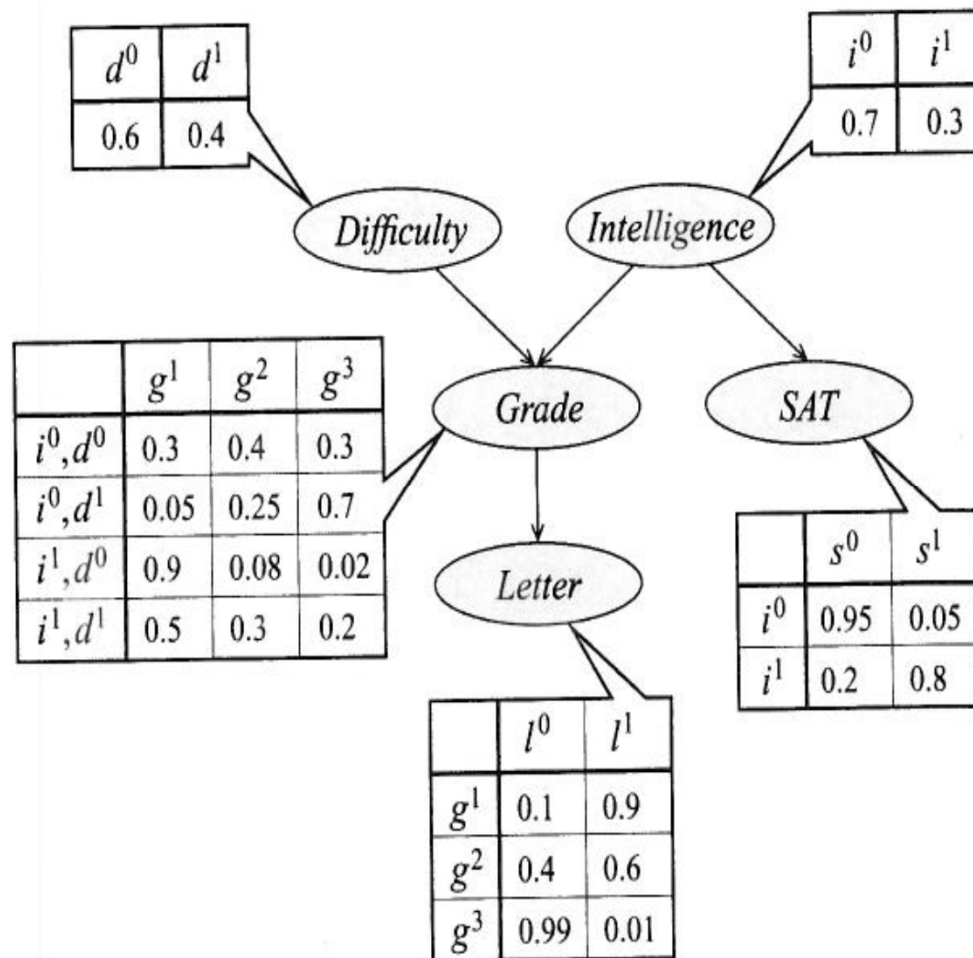
Markov Chain Monte Carlo (MCMC) Methods

- Forward sampling and likelihood weighted sampling methods are defined for BNs only.
 - Also, when evidence is mostly at the leaf nodes, the methods largely sample from the prior distribution
- Another approach is to generate a sequence (chain)
 - Called Markov chain because the next element depends only on the current sample
 - Monte Carlo because of random sampling
- Applies to both directed and undirected models
- Has many practical applications

Gibbs Sampling (Gibbs Chain)

- Reduce all factors according to evidence
- Generate an initial sample
 - Say from forward sampling; reduce factors if evidence variables are given
- Iterate over unobserved variables, one at a time
 - Sample a new value, given the current sample values for all other variables
 - Consider variable X_i ; sample from $P_{\Phi}(X_i | \mathbf{x}_{-i})$; \mathbf{x}_{-i} means assignment of all variables in X except variable X_i .
 - Example and algorithm to follow
- At each step, we are sampling from a posterior distribution (though not necessarily the correct one)
- Can be shown that the process converges to the actual posterior distribution

Example 12.4



Example 12.4

- Let evidence be s^1 and l^0
- Reduced factors are $P(I)$, $P(D)$, $P(G|I,D)$, $P(s^1 | I)$, $P(l^0|G)$
- Let first sample, at time 0, by forward sampling, be d^1, i^0, g^2
- Now sample in order of G, I, D (note: need not start from root)
- Sample for G is drawn from $P_{\Phi}(G | d^1, i^0, s^1, l^0)$

$$\begin{aligned} P_{\Phi}(G | d^1, i^0) &= \frac{P(i^0)P(d^1)P(G | i^0, d^1)P(l^0 | G)P(s^1 | i^0)}{\sum_g P(i^0)P(d^1)P(g | i^0, d^1)P(l^0 | g)P(s^1 | i^0)} \\ &= \frac{P(G | i^0, d^1)P(l^0 | G)}{\sum_g P(g | i^0, d^1)P(l^0 | g)}. \end{aligned}$$

- Multiply all factors containing G and normalize
- Sample from this distribution, let the sample yield $G = g^3$
- Now sample I from $P_{\Phi}(I | d^1, g^3)$;
- Sample D next and keep iterating
- Distribution from which samples are drawn, get closer and closer to the posterior distribution

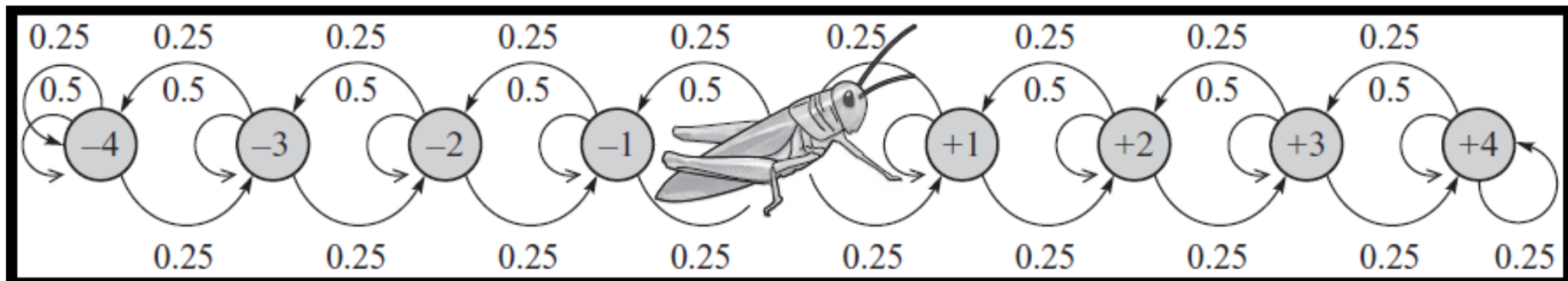
Algorithm 12.4 Generating a Gibbs chain trajectory

Procedure Gibbs-Sample (
 X // Set of variables to be sampled
 Φ // Set of factors defining P_Φ
 $P^{(0)}(X)$, // Initial state distribution
 T // Number of time steps
)

- 1 Sample $x^{(0)}$ from $P^{(0)}(X)$
- 2 **for** $t = 1, \dots, T$
- 3 $x^{(t)} \leftarrow x^{(t-1)}$
- 4 **for** each $X_i \in X$
- 5 Sample $x_i^{(t)}$ from $P_\Phi(X_i \mid x_{-i})$
- 6 // Change X_i in $x^{(t)}$
- 7 **return** $x^{(0)}, \dots, x^{(T)}$

Markov Chains

- Defined by a transition function $T(\mathbf{x} \rightarrow \mathbf{x}')$ between a pair of states $(\mathbf{x}, \mathbf{x}')$ which defines the probability of going from current state \mathbf{x} to new state \mathbf{x}' . A state is given by assignments to variables.
 - Note T will have n^2 entries if \mathbf{X} can take n values
 - Can be viewed as a matrix
- Homogeneous Markov Chain
 - Transition probability does not change over time
- Grasshopper Example
 - State: 9 integers from -4 to +4
 - Initial position: 0
 - At each instance, $T(i \rightarrow i) = .5$, $T(i \rightarrow i-1) = .25$, $T(i \rightarrow i+1) = .25$
 - At two ends, can not jump beyond (stays in the same state)
 - $T(4 \rightarrow 4) = .75$
 - Write as a transition matrix



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

At $t=0$, $P(X^0=0) = 1$

At $t = 1$, $P(X^1 = 0) = .5$, $P(X^1 = 1) = .25$, $P(X^1 = -1) = .5$

At $t = 2$, $P(X^2 = 0) = .5 \times .5 + .25 \times .25 + .25 \times .25 = .375$

$P(X^2 = 1 \text{ or } -1) = .5 \times .25 + .25 \times .5 = .25$

$P(X^2 = 2 \text{ or } -2) = .25 \times .25 = .0625$

Position probability converges to a nearly uniform distribution
with time for this example

MCMC Sampling

- Generate a chain by sampling from the distribution
 - Sample $x^{(t)}$ from distribution $P^{(t)}$
 - Does $P^{(t)}$ converge and if so, to the desired distribution

Algorithm 12.5 Generating a Markov chain trajectory

Procedure MCMC-Sample (
 $P^{(0)}(X)$, // Initial state distribution
 T , // Markov chain transition model
 T // Number of time steps
)

- 1 Sample $x^{(0)}$ from $P^{(0)}(X)$
- 2 **for** $t = 1, \dots, T$
- 3 Sample $x^{(t)}$ from $T(x^{(t-1)} \rightarrow X)$
- 4 **return** $x^{(0)}, \dots, x^{(T)}$

Stationary Distribution

- At convergence, we expect:

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

- Stationary Distribution

A distribution $\pi(X)$ is a stationary distribution for a Markov chain T if it satisfies:

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

*A stationary distribution is also called an *invariant distribution*.*

- In linear algebra formulation: $T \pi(x) = \pi(x)$; *i.e.* the stationary distribution is an eigenvector of the transition matrix with eigenvalue = 1

Next Class

- Read sections 12.3, 6.2, 15.1 of the KF book