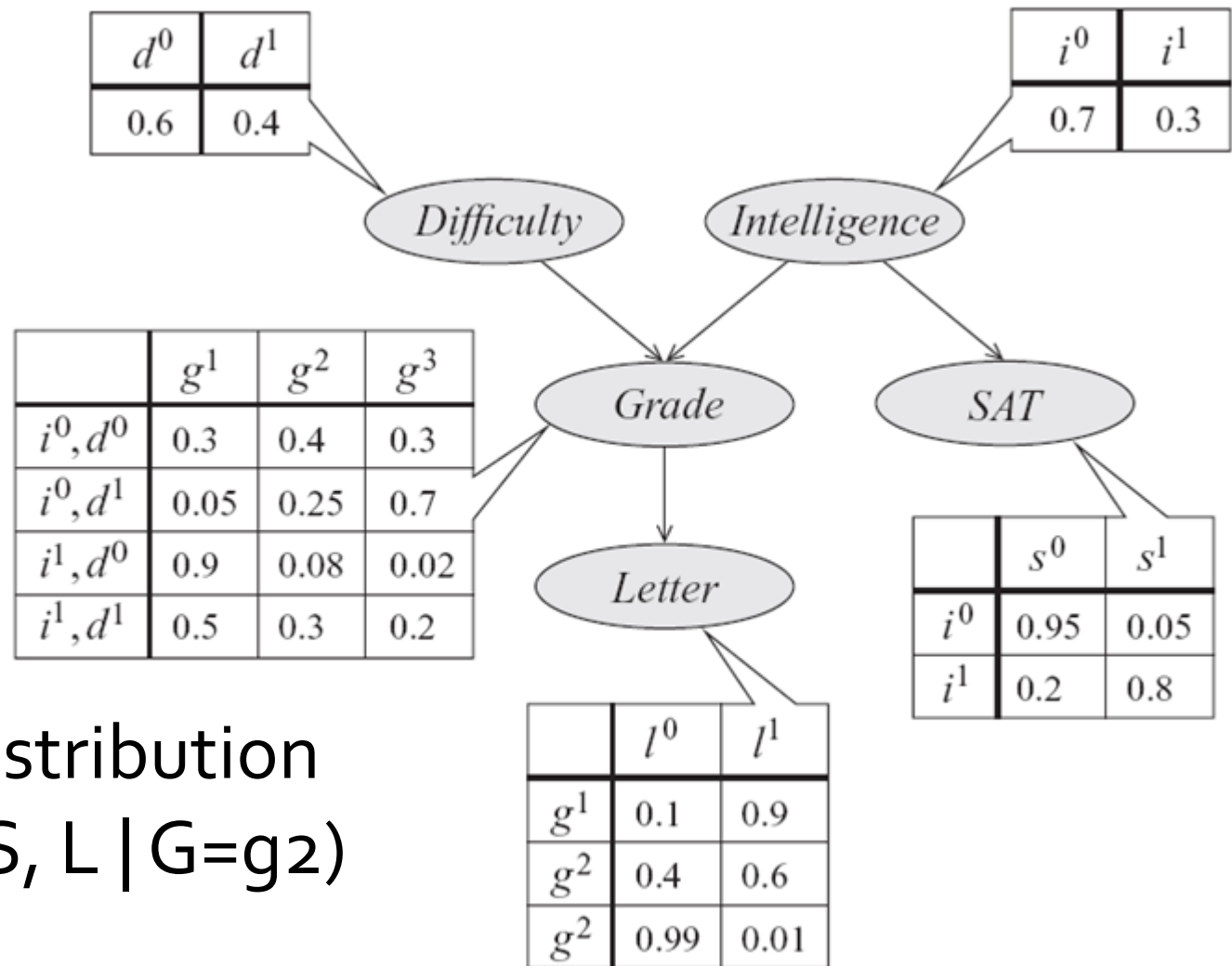


MCMC based Sampling

Our Goal (reminder)

- We need independent samples to estimate a desired distribution (usually posterior distribution, $p(Y|e)$)
- We can setup a Markov chain that converges to a stationary distribution
- Satisfying detailed balance is an easy way to guarantee convergence to equilibrium

Example



- Desired distribution $P(D, I, S, L | G=g^2)$

Problem Setup

- Usually the state space is huge but in our toy example:

$$|\mathbb{I}| = |D| \times |I| \times |S| \times |L| \times |G| = 2 \times 2 \times 2 \times 2 \times 3 = 48$$

- Samples are shown as $\mathbf{x}:(d,i,s,l,g)$
- Given our graphical model we can simply evaluate every sample as:

$$p(\mathbf{x}) = p(d)p(i)p(s|i)p(g|d,i)p(l|g)$$

Average of Samples Converge to the Expectation

- Ergodicity (special case of law of large numbers). If a Markov process is positive recurrent with invariant distribution π then

$$P\left(\frac{1}{n} \sum_{k=0}^{n-1} f(X_k) \rightarrow \bar{f} \text{ as } n \rightarrow \infty\right) = 1$$

- Where

$$\bar{f} = E_{\pi}(f) = \sum_{i \in I} \pi_i f_i$$

- And π is the unique invariant distribution

Proposal Distribution (reminder)

- We cannot always sample efficiently from $P(X)$
- But we might be able to evaluate $P(X)$ efficiently
- In that case, we could sample efficiently from some other “simpler” distribution called the *proposal distribution* $Q(X)$

Proposal Distribution (cont.)

- if $Q(x) \neq 0$ where $P(x) \neq 0$

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

How to find desired distribution?

- Several ways we can do this with MCMC
 - Metropolis
 - Metropolis Hasting
 - Gibbs Sampling

General form of MCMC

1. Sample a point from a *proposal distribution*
 $q(y | x)$

2. Compute the *importance ratio*

$$r = \frac{p(y) q(x | y)}{p(x) q(y | x)}$$

3. Move to the new state with an *transition probability* (related to importance ratio)

$$P(x \rightarrow y) = q(y | x) \{r \wedge 1\}$$

Metropolis Algorithm

- The probability of moving from one state to another *must be symmetric*:

$$q(\mathbf{x} | \mathbf{y}) = q(\mathbf{y} | \mathbf{x})$$

Metropolis Algorithm

- Importance ratio

$$r = \frac{p(y) q(x | y)}{p(x) q(y | x)} = \frac{p(y)}{p(x)}$$

- Transition probability

$$P(x \rightarrow y) = q(y | x) \{r \wedge 1\}$$

Metropolis Algorithm

```
X ← randomValue ()
while (1) :
    Y = generateSample (q (Y | X) )
    r = p (Y) / p (X)
    if (r > 1) :
        X = Y
    else:
        t = generateSample (uniform (0, 1) )
        if (t < r) :
            X = Y
```

Convergence of Metropolis Algorithm

$$p(x)P(x \rightarrow y) = p(x)q(y/x) \left\{ \frac{p(y)}{p(x)} \wedge 1 \right\}$$

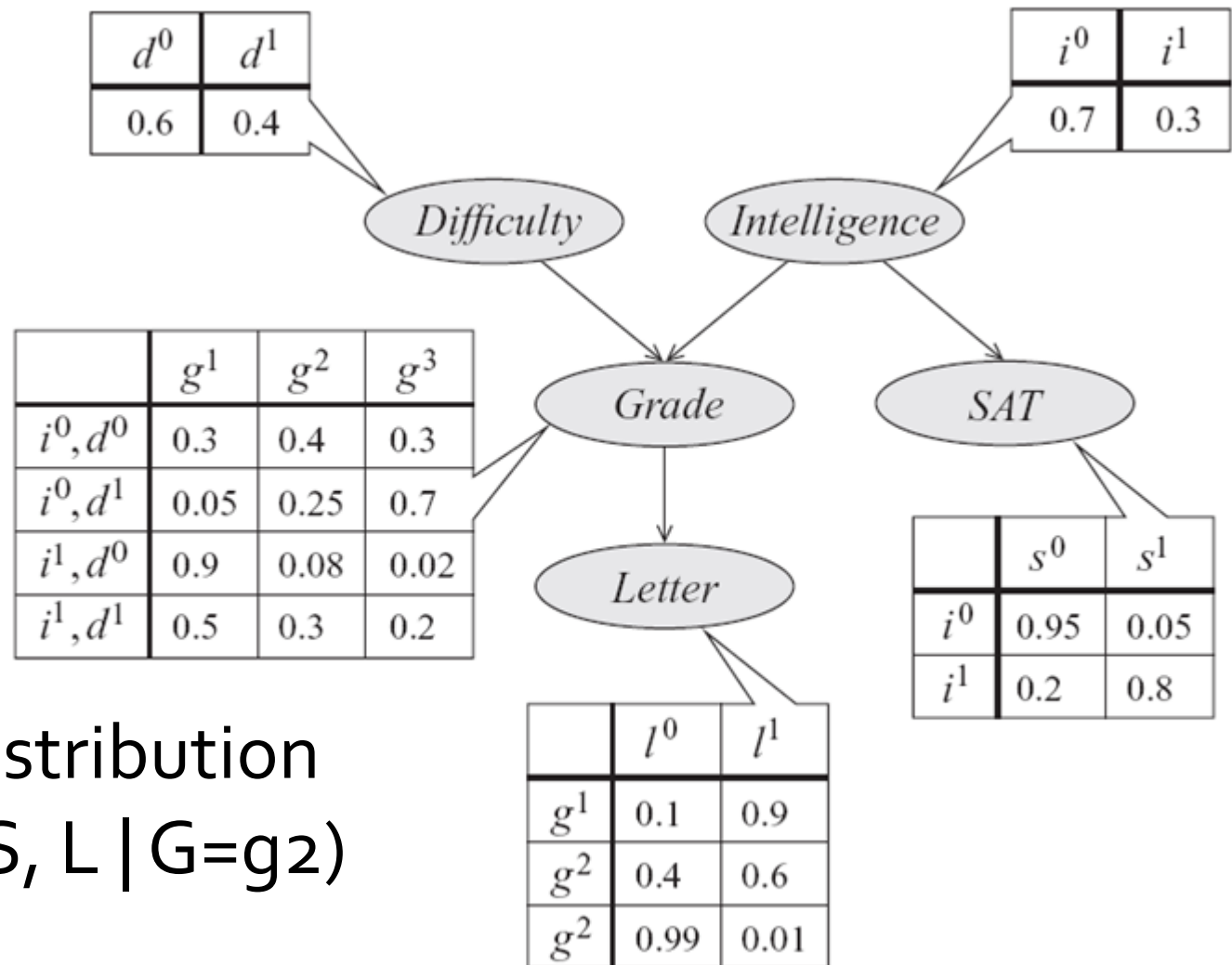
assume $p(y) > p(x)$ then

$$= p(x)q(y/x) \times 1 = p(x)q(x/y)$$

$$= p(x)q(x/y) \frac{p(y)}{p(y)} = p(y)q(x/y) \frac{p(x)}{p(y)}$$

$$= p(y)q(x/y) \left\{ \frac{p(x)}{p(y)} \wedge 1 \right\} = p(y)P(y \rightarrow x)$$

Example



- Desired distribution $P(D, I, S, L \mid G=g^2)$

Problem Setup

- Usually the state space is huge but in our toy example:

$$|\mathbb{I}| = |D| \times |I| \times |S| \times |L| \times |G| = 2 \times 2 \times 2 \times 2 \times 3 = 48$$

- Samples are shown as $\mathbf{x}:(d,i,s,l,g)$
- Given our graphical model we can simply evaluate every sample as:

$$p(\mathbf{x}) = p(d)p(i)p(s|i)p(g|d,i)p(l|g)$$

Metropolis(example)

- Let proposal distribution be uniform

- Start with random $\mathbf{x}:(d^0, i^1, s^1, l^1, g^2)$

$$p(x) = 0.6 \times 0.3 \times 0.8 \times 0.08 \times 0.6 = 0.007$$

- Obtain \mathbf{y} by uniformly sampling from I

$$\mathbf{y}:(d^1, i^1, s^0, l^0, g^2)$$

$$p(y) = 0.4 \times 0.3 \times 0.2 \times 0.3 \times 0.4 = 0.003$$

$$\frac{p(y)}{p(x)} \wedge 1 = \frac{0.003}{0.007} \wedge 1 = 0.42$$

- Draw a random value between $[0,1]$. If it is " $>$ " than 0.42 reject it. Let say 0.32, so accept \mathbf{y} .

Metropolis(example)

- Again obtain the next \mathbf{y} by uniformly sampling from I .

$$\mathbf{y}:(d^0, i^0, s^0, l^0, g^2)$$

$$p(\mathbf{y}) = 0.6 \times 0.7 \times 0.95 \times 0.4 \times 0.4 = 0.064$$

$$\frac{p(\mathbf{y})}{p(\mathbf{x})} \wedge 1 = \frac{0.064}{0.003} \wedge 1 = 1$$

- So accept \mathbf{y} and so on

Metropolis-Hastings Algorithm

- The proposal distribution need not be symmetric
- Now the proposal distribution is factored into the importance ratio
- Follows the general form introduced earlier
- This is more general (and useful) than Metropolis algorithm

Metropolis-Hastings

- Importance ratio

$$r = \frac{p(y)q(x/y)}{p(x)q(y/x)}$$

- Transition probability

$$P(x \rightarrow y) = q(y/x) \{ r \wedge 1 \}$$

Metropolis-Hastings Algorithm

```
X ← randomValue ()
while (1) :
    Y = generateSample (q (Y | X) )
    r = p (Y) q (X | Y) / p (X) q (Y | X)
    if (r > 1) :
        X = Y
    else:
        t = generateSample (uniform (0, 1) )
        if (t < r) :
            X = Y
```

Convergence of Metropolis-Hasting Algorithm

$$p(x)P(x \rightarrow y) = p(x)q(y|x) \left\{ \frac{p(y)q(x|y)}{p(x)q(y|x)} \wedge 1 \right\}$$

$$\text{assum } p(y)q(x|y) > p(x)q(y|x)$$

$$= p(x)q(y|x) = p(x)q(y|x) \frac{p(y)q(x|y)}{p(y)q(x|y)}$$

$$= p(y)q(x|y) \frac{p(x)q(y|x)}{p(y)q(x|y)} = p(y)q(x|y) \frac{p(x)q(y|x)}{p(y)q(x|y)}$$

$$= p(y)q(x|y) \left\{ \frac{p(x)q(y|x)}{p(y)q(x|y)} \wedge 1 \right\}$$

$$= p(y)P(y \rightarrow x)$$

Gibbs Sampling Algorithm

- Special case of Metropolis-Hastings algorithm
- The proposal distribution has a given form (i.e. it is not designed on a problem by problem basis)
- Samples the components of the outcome vector one at a time using the marginal distribution, where all other components are fixed to values from previous samples

Proposal Distribution of Gibbs Sampling

- Proposal distribution

$$q(y/x) = \begin{cases} p(y_j/x_{-j}) & y_{-j} = x_{-j} \\ 0 & \textit{otherwise} \end{cases}$$

- Importance ratio is unity. We *always* accept.

$$r = \frac{p(y)q(x|y)}{p(x)q(y|x)} = 1$$

Gibbs Sampling Algorithm

```
X←randomValue()  
while(1):  
    for(j=0;j<len(X);j++)  
        y=generateSample(p(y|X[0:j],X[j+1:len(x)]))  
        X[j]=y
```


Proof that acceptance rate is unity

$$\begin{aligned} r &= \frac{p(y)q(x/y)}{p(x)q(y/x)} \\ &= \frac{p(y)p(x_j/y_{-j})}{p(x)p(y_j/x_{-j})} = \frac{p(y)p(x_j/x_{-j})}{p(x)p(y_j/y_{-j})} \\ &= \frac{p(y)p(x_j, x_{-j})p(y_{-j})}{p(x)p(y_j, y_{-j})p(x_{-j})} = \frac{p(y)p(x)p(y_{-j})}{p(x)p(y)p(x_{-j})} \\ &= \frac{p(y_{-j})}{p(x_{-j})} = 1 \end{aligned}$$

Gibbs Sampling

- Let start with $\mathbf{x}:(d^0, i^1, s^1, l^1, g^2)$
- We will sample D, I, S, L and G in a round robin manner
- Sample D:
 - $s_1: p(d^0, i^1, s^1, l^1, g^2) = 0.6 \times 0.3 \times 0.8 \times 0.08 \times 0.6 = 0.007$
 - $s_2: p(d^1, i^1, s^1, l^1, g^2) = 0.4 \times 0.3 \times 0.8 \times 0.3 \times 0.6 = 0.017$
 - Generate a random number between $[0, 0.007 + 0.017]$, if it is greater than 0.007 move to s_2 otherwise move to s_1 . let say random value is 0.011. so current state will become $\mathbf{x}:(d^1, i^1, s^1, l^1, g^2)$

Gibbs Sampling

- Sample 1:
 - $s_1: p(d^1, i^0, s^1, l^1, g^2) = 0.4 \times 0.7 \times 0.05 \times 0.25 \times 0.6 = 0.002$
 - $s_2: p(d^1, i^1, s^1, l^1, g^2) = 0.4 \times 0.3 \times 0.8 \times 0.3 \times 0.6 = 0.017$
 - Generate a random number between $[0, 0.002 + 0.017]$, if it is greater than 0.002 move to s_2 otherwise move to s_1 . let say random value is 0.001. so current state will become $\mathbf{x}: (d^1, i^0, s^1, l^1, g^2)$
- **Sample S and so on**

Practical issues

- Learning how a tool works is one thing, using it in a practical situation is another.
- MCMC is no different.

Convergence

- Convergence and ergodicity theorems state

$$P(X_n = j) \rightarrow \pi_j \text{ as } n \rightarrow \infty$$

$$\frac{1}{n} \sum_{k=0}^{n-1} f(X_k) \rightarrow \bar{f} \text{ as } n \rightarrow \infty$$

- This is nice, but they don't say anything about *how fast* they converge

Convergence

- How can we tell if our sample has accurately characterized our desired distribution?
- How big should n be before we trust our result?
- How long does this MCMC thing going to take?

Mixing time

- Mixing is a measure of how long a Markov process takes to get near its equilibrium.
- There are many analytic ways to calculate this, but only for completely characterized Markov chains.
- For sampling, we would like our process to converge quickly from an arbitrary point in the space. This is called “mixing well.”

Burn In

- We start the Markov chain from a random point in the sample space.
- This point and the points in its neighborhood might be very unlikely according to our distribution.
- Run the Markov chain for many iterations before using the sampled points. This is called *burn in*.

Burn In

- How long should our burn in last?
 1. Analytically evaluate the convergence rate of our chain
 1. Usually results in overly pessimistic estimates
 2. Use convergence diagnostics
 1. Do not guarantee convergence
 3. Use perfect simulation
 1. Only valid for specific types of problems

Thinning

- Consecutive outputs from the chain can be highly correlated
- Saving all of the sampled points can be expensive
- We can save only every k outputs, which is called *thinning*
- Sample of our samples

Review

- Markov chains are ergodic

$$P\left(\frac{1}{n} \sum_{k=0}^{n-1} f(X_k) \rightarrow \bar{f} \text{ as } n \rightarrow \infty\right) = 1$$

- MCMC is a technique for importance sampling
- Metropolis requires a symmetric proposal distribution
- Metropolis-Hasting does not

Review

- Gibbs sampling is special case of Metropolis Hastings that always accepts the proposal
- In implementing MCMC, convergence and independence are major concerns
 - Burn in (convergence rate, convergence diagnostics , perfect sampling)
 - Thinning