# Monte Carlo Markov Chain 1

# MCMC

Limitations of LW:

- <u>Evidence</u> affects sampling only for nodes that are its descendants
- For nondescendants, the <u>weights</u> account for the effect of the evidence
- If evidence is at the leaves, we are sampling from the prior distribution (and not the posterior which is what we want)

# MCMC

Strategy used by MCMC

- Generate a <u>sequence</u> of samples
- Initial samples generated from the prior
- Successive samples generated progressively closer to the posterior

Applies to both directed and undirected models. We'll use a distribution $P_\Phi$ defined in terms of a set of factors $\Phi$
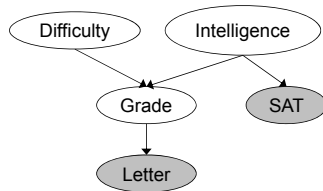
# Gibbs Sampling

# Gibbs Sampling

Example: Suppose we have as evidence *SAT = High and Letter = Weak* (nodes are shaded grey)



Factors:
- P(I)
- P(D)
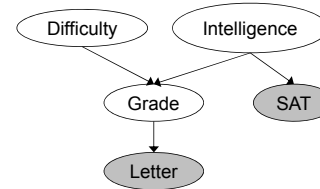- P(G | I,D)

Reduced Factors:
- P(S=high | I)
- P(L=weak | G)

Eliminate all rows that are inconsistent with the evidence in all factors (see pg 111 of textbook)
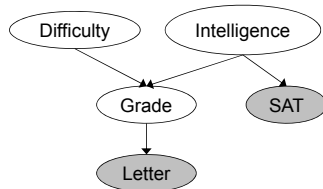
5

---

# Gibbs Sampling



Start with an initial sample eg: $x^{(0)}$ = *(D = high, I = low, G = B, S = high, L = weak)*

- *D*, *I* and *G* could be set in any way, for instance by forward sampling, to get $D^{(0)}$ = *high*, $I^{(0)}$ = *low*, $G^{(0)}$ = *B*

- *S=high* and *L=weak* are observed

6

---

# Gibbs Sampling



Resample non-evidence nodes, one at a time, in some order eg. G, I, D.

If we sample $X_i$, keep other nodes clamped at the values of the current state *(D = high, I = low, G = B, S = high, L = weak)*
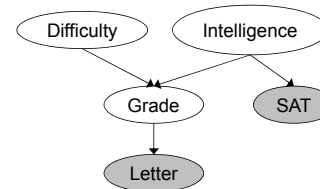
To sample $G^{(1)}$, we compute $P_\Phi$(G | D=high, I=low, S=high, L=weak):

$$P_\Phi(G|D = high, I = low, S = high, L = weak)$$
$$= \frac{P(I = high)P(D = low)P(G|I = low, D = High)P(L = low|G)P(S = high|I = low)}{\sum_G P(I = high)P(D = low)P(G|I = low, D = High)P(L = low|G)P(S = high|I = low)}$$
$$= \frac{P(G|I = low, D = high)P(Letter = weak|G)}{\sum_G P(G|I = low, D = high)P(Letter = weak|G)}$$

7

---

# Gibbs Sampling



- Suppose we obtain $G^{(1)}$ = C.
- Now sample $I^{(1)}$ from $P_\Phi$(I | D=high, G=C, S=high, L=weak). Note it is conditioned on $G^{(1)}$=C
- Say we get $I^{(1)}$=high
- Now sample $D^{(1)}$ from $P_\Phi$(D | G=C, I = high, S=high, L=weak). Say you get $D^{(1)}$ = high
- The first iteration of sampling produces $x^{(1)}$ = (I = high, D = high, G = C, S=high, L=weak)
- Iterate...

8

## Gibbs Sampling

- $P_\Phi(G \mid D=high, I=low, S=high, L=weak)$ takes downstream evidence *L=weak* into account (makes it closer to the posterior distribution $P(X \mid \textbf{\textit{e}})$)
- Early on, $P_\Phi(G \mid D=high, I=low, S=high, L=weak)$ very much like the prior *P(X)* because it uses values for *I* and *D* sampled from *P(X)*
- On next iteration, resampling *I* and *D* conditioned on new value of *G* brings the sampling distribution closer to the posterior
- <span style="color:red">Sampling distribution gets progressively closer and closer to the posterior</span>

9

---

## Gibbs Sampling

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

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

10

---

## Gibbs Sampling

Gibbs sampling with evidence
- Reduce all factors by the observations **e**
- The distribution $P_\Phi$ corresponds to *P(X|**e**)*

11

---

## Markov Chains

12

# Markov Chains

- (Informally) A Markov chain is a graph of states over which the sampling algorithm takes a random walk
- Note: the graph is not the graphical model but a graph over the possible assignments to a set of variables **X**

13

---

# Markov Chains

- A Markov chain is defined via a state space *Val(X)* and a model that defines, for every state $x \in Val(X)$ a next-state distribution over *Val(X)*.
- More precisely, the transition model $\mathcal{T}$ specifies for each pair of states $x$, $x'$ the probability $\mathcal{T}(x \to x')$ of going from $x$ to $x'$.
- A homogeneous Markov chain is one where the system dynamics do not change over time

14

---
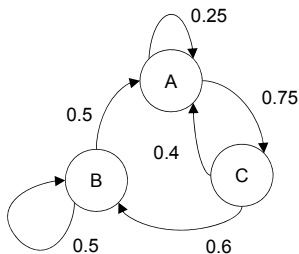
# Markov Chains

Example of a Markov Chain with Val(X)={A,B,C}:

State Transition Diagram View



Conditional Probability Distribution View

| $X_{t-1}$ | $X_t$ | $P(X_t|X_{t-1})$ |
|---|---|---|
| A | A | 0.25 |
| A | B | 0 |
| A | C | 0.75 |
| B | A | 0.5 |
| B | B | 0.5 |
| B | C | 0 |
| C | A | 0.4 |
| C | B | 0.6 |
| C | C | 0 |

15

---

# Markov Chains

- Random sampling process defines a random sequence of states $\mathbf{x}^{(0)}$, $\mathbf{x}^{(1)}$, $\mathbf{x}^{(2)}$, …
- $\mathbf{X}^{(t)}$ is a random variable:
- Need initial state distribution $P^{(0)}(\mathbf{X}^{(0)})$
- Probability that next state is $\mathbf{x}'$ can be computed as:

$$P^{(t+1)}(\mathbf{X}^{(t+1)} = \mathbf{x}') = \sum_{x \in Val(X)} P^{(t)}(\mathbf{X}^{(t)} = \mathbf{x})\mathcal{T}(\mathbf{x} \to \mathbf{x}')$$

Sum over all states that the chain could have been at time t

Probability of transition from x to x'

16

---

13

14

15

16

4

# Markov Chains

How to generate a Markov Change Monte Carlo trajectory:

> **Procedure** MCMC-Sample (
>     $P^{(0)}(\mathbf{X})$, // Initial state distribution
>     $\mathcal{T}$,     // Markov chain transition model
>     T     // Number of time steps
> )
> 1.   Sample $x^{(0)}$ from $P^{(0)}(\mathbf{X})$
> **2.**   **for** t = 1, …, T
> 3.   Sample $\mathbf{x}^{(t)}$ from $\mathcal{T}(\mathbf{x}^{(t-1)} \rightarrow \mathbf{X})$
> **4.**   **return** $x^{(0)}, \ldots, x^{(T)}$

The big question: does $P^{(t)}$ converge and what to?

17

---

# Markov Chains

- When the process converges, we expect:

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

- A distribution $\pi(\mathbf{X})$ is a stationary distribution for a Markov chain $\mathcal{T}$ if it satisfies:

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

- A stationary distribution is also called an invariant distribution

18

---

# Markov Chains

Another example:



0.25
0.7
0.75
0.5
0.5
0.3

$X_1$  $X_2$  $X_3$

To find the stationary distribution:

$\pi(x_1) = 0.25\pi(x_1)+0.5\pi(x_3)$

$\pi(x_2) = 0.7\pi(x_2)+0.5\pi(x_3)$

$\pi(x_3) = 0.75\pi(x_1)+0.3\pi(x_2)$

$\pi(x_1) + \pi(x_2) + \pi(x_3) = 1$

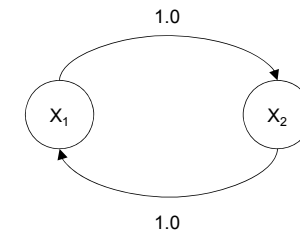Solving these simultaneous equations gives: $\pi(x_1) = 0.2$, $\pi(x_2) = 0.5$, $\pi(x_3) = 0.3$

19

---

# Markov Chains

- Bad news: no guarantee that MCMC sampling process converges to a stationary distribution
- Example of a periodic Markov chain (periodic = fixed cyclic behavior)
  - Start with $P^{(0)}(x_1) = 1$
  - $P^{(t)}(x_1) = 1$ if t is even
  - $P^{(t)}(x_2) = 1$ if t is odd

1.0
1.0

$X_1$  $X_2$

20

## Markov Chains

- No guarantee that stationary distribution is unique – depends on $P^{(0)}$
  - This happens if the chain is reducible: has states that are not reachable from each other
- We will restrict our attention to Markov chains that have a stationary distribution which is reached from any starting distribution $P^{(0)}$

21

## Markov Chains

- To meet this restriction, we need the chain to be regular
- A Markov chain is said to be regular if there exists some number $k$ such that, for every $x, x' \in Val(\textbf{X})$, the probability of getting from $x$ to $x'$ in exactly $k$ steps is > 0
- Theorem 12.3: If a finite state Markov chain $\mathcal{T}$ is regular, then it has a unique stationary distribution

22

## Markov Chains

- Define $\mathcal{T}_i$ to be a transition model called a kernel
- For graphical models, define a kernel $\mathcal{T}_i$ for each variable $X_i \in \textbf{X}$
- Define $\textbf{X}_{-i} = \mathcal{X} - \{X_i\}$ and let $\textbf{x}_i$ denote an instantiation to $\textbf{X}_i$
- The model $\mathcal{T}_i$ takes a state $(\textbf{x}_{-i}, x_i)$ and transitions to a state $(\textbf{x}_{-i}, x_i')$

23

## Gibbs Sampling Revisited

24

## Gibbs Sampling Revisited

How do we use MCMC on a graphical model?
- Want to generate samples from the posterior $P(\boldsymbol{X}|\boldsymbol{E}=\boldsymbol{e})$ where $\boldsymbol{X}=\mathcal{X}-\boldsymbol{E}$
- Define a chain where $P(\boldsymbol{X}|\boldsymbol{e})$ is the stationary distribution
- States are instantiations $\boldsymbol{x}$ to $\mathcal{X}-\boldsymbol{E}$
- Need transition function that converges to stationary distribution $P(\boldsymbol{X}|\boldsymbol{e})$
- For convenience: define $P_\Phi = P(\boldsymbol{X}|\boldsymbol{e})$ where the factors in $\Phi$ are reduced by the evidence $\boldsymbol{e}$

25

## Gibbs Sampling Revisited

Using the MCMC framework, the transition model for Gibbs Sampling is:

$$\mathcal{T}_i((\boldsymbol{x}_{-i}, x_i) \to (\boldsymbol{x}_{-i}, x_i^{'})) = P(x_i^{'} \mid \boldsymbol{x}_{-i})$$

And the posterior distribution $P_\Phi(\boldsymbol{X}) = P(\mathcal{X}|\boldsymbol{e})$ is a stationary distribution of this process

26

## Gibbs Sampling Revisited

Gibbs sampling on a Bayesian network is efficient
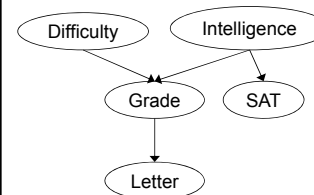Note: $Pa(x_i)$ = Parents of $x_i$, $Ch(x_i)$ = Children of $x_i$

$P(X_i|x_1, \dots, x_{i-1}, x_{i+1}, x_n)$
$= \dfrac{P(x_1, \dots, x_{i-1}, x_i, x_{i+1}, x_n)}{P(x_1, \dots, x_{i-1}, x_{i+1}, x_n)} = \dfrac{P(x_1, \dots, x_{i-1}, x_i, x_{i+1}, x_n)}{\sum_{x_i} P(x_1, \dots, x_{i-1}, x_i, x_{i+1}, x_n)}$

$= \dfrac{\prod_{j=1}^n P(x_j|Pa(x_j))}{\sum_{x_i} \prod_{j=1}^n P(x_j|Pa(x_j))}$

$= \dfrac{\prod_{x_j \notin \{x_i, Ch(x_i)\}} P(x_j|Pa(x_j)) P(x_i|Pa(x_i)) \prod_{x_k \in Ch(x_i)} P(x_k|x_i, other\ parents)}{\prod_{x_j \notin \{x_i, Ch(x_i)\}} P(x_j|Pa(x_j)) \sum_{x_i} P(x_i|Pa(x_i)) \prod_{x_k \in Ch(x_i)} P(x_k|x_i, other\ parents)}$

$= \dfrac{P(x_i|Pa(x_i)) \prod_{x_k \in Ch(x_i)} P(x_k|x_i, other\ parents)}{\sum_{x_i} P(x_i|Pa(x_i)) \prod_{x_k \in Ch(x_i)} P(x_k|x_i, other\ parents)}$

Depends only on the
CPDs of $X_i$ and its children

27

## Gibbs Sampling Revisited

Difficulty   Intelligence

Grade   SAT

Letter

Student Example revisited:

Define:

$\mathcal{T}((I,G,D,S=high,L=weak) \to (I', G, D, S=high, L=weak)) = P(I|G,D,S=high,L=weak)$

Sample from the distribution below:

$P(I' | G, D, S = high, L = weak)$
$= \dfrac{P(I')P(G \mid I', D)P(S = high \mid I')}{\sum_{I''} P(I''= i'')P(G \mid I''= i'', D)P(S = high \mid I''= i'')}$

28

## Gibbs Sampling Revisited

Block Gibbs Sampling
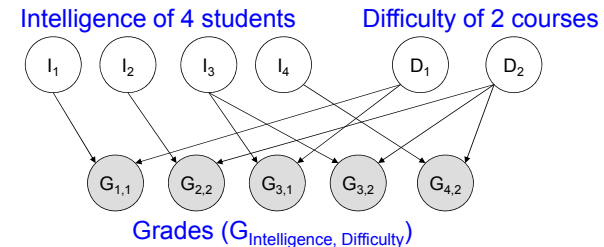
- Can sample more than a single variable $X_i$ at a time
- Partition $X$ into disjoint blocks of variables $X_1, ..., X_k$
- Then sample $P_\Phi(X_i \mid X_1=x_1, ..., X_{i-1}=x_{i-1}, X_{i+1}=x_{i+1}, ..., X_k=x_k)$
- Takes longer range transitions

29

29

## Gibbs Sampling Revisited

Intelligence of 4 students    Difficulty of 2 courses



Grades ($G_{Intelligence, Difficulty}$)

- Step t: Sample all of the I variables as a block, given Ds and Gs (since Is are conditionally independent from each other given Ds)
- Step t+1: Sample all of the D variables as a block, given Is and Gs (since Ds are conditionally independent of each other given Is)    30

30

## Gibbs Sampling Revisited

Need to compute $P_\Phi(X_i \mid X_1=x_1, ..., X_{i-1}=x_{i-1}, X_{i+1}=x_{i+1}, X_k = x_k)$

- Efficient if variables in each block (eg. **I**) are independent given the variables outside the block (eg. **D**)
- In general, full independence is not essential – need some sort of structure to the block-conditional distribution

31

31

## Gibbs Sampling Revisited

- Gibbs chain not necessarily regular and may not converge to a unique stationary distribution
- Only guaranteed to be regular if $P(X_i \mid X_{-i})$ is positive for every value of $X_i$
- Theorem 12.4: Let $\mathcal{H}$ be a Markov network such that all of the clique potentials are strictly positive. Then the Gibbs-sampling Markov chain is regular.

32

32