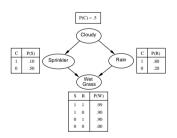
Lecture 12: Particle-based inference: Gibbs sampling

- Gibbs sampling
- Markov chains
- Markov Chain Monte Carlo (MCMC) methods

January 29, 2007 1 COMP-526 Lecture 12

A different idea



- Suppose we want to compute P(R|S=1)
- We generate <u>one sample</u>, with the given evidence variables instantiated correctly
- Then we keep changing it!
- If we are careful, we will get samples from the correct distribution

Recall from last time: Particle-based inference

- ullet Suppose we have evidence E=e and we want to know p(Y|E=e) for some query variables Y
- Particle-based methods will generate particles and then compute sufficient statistics to estimate this answer
- Likelihood weighting has an easy way of producing samples: go through the Bayes net in the direction of the arcs, sample nodes without evidence and set the value for evidence variables
- ullet Since these samples are from a "mutilated" Bayes net, NOT from p(Y|E=e) each particle must have a weight. The weights are used instead of counts in the probability estimation.
- But these weights can get very small, and then we would need to sample a lot of data to get good estimates.

January 29, 2007

2

COMP-526 Lecture 12

Gibbs sampling

- 1. Initialization
 - ullet Set evidence variables Ej, to the observed values e
 - Set all other variables to random values (e.g. by forward sampling, uniform sampling...)

This gives us a sample x_1, \ldots, x_n .

- 2. Repeat (as much as wanted)
 - Pick a non-evidence variable X_i uniformly randomly)
 - Sample x_i' from $p(X_i|x_1,\ldots,x_{i-1},x_{i+1},\ldots,x_n)$.
 - Keep all other values: $x'_i = x_i, \forall j \neq i$
 - The new sample is x'_1, \ldots, x'_n
- Alternatively, you can march through the variables in some predefined order

4

January 29, 2007 3 COMP-526 Lecture 12

January 29, 2007

COMP-526 Lecture 12

Why Gibbs works in Bayes nets

- The key step is sampling according to $p(X_i|x_1,\dots,x_{i-1},x_{i+1},\dots,x_n).$ How do we compute this?
- In Bayes nets, we know that a variable is conditionally independent of all other *given its Markov blanket* (parents, children, spouses)

$$p(X_i|x_1,\ldots,x_{i-1},x_{i+1},\ldots,x_n)=p(X_i|\mathsf{MarkovBlanket}(X_i))$$

- So we need to sample from $P(X_i|\mathsf{MarkovBlanket}(X_i))$
- Let $Y_i, j = 1, \ldots, k$ be the children of X_i . We can show that:

$$p(X_i=x_i|\mathsf{MarkovBlanket}(X_i)) \quad \propto \quad p(X_i=x_i|\mathsf{Parents}(X_i)) \cdot \\ \cdot \quad \prod_{j=1}^k p(Y_j=y_j|\mathsf{Parents}(Y_j))$$

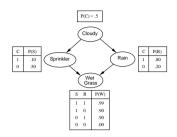
January 29, 2007 5 COMP-526 Lecture 12

Analyzing Gibbs sampling

- Consider the variables X_1, \ldots, X_n . Each possible assignment of values to these variables is a state of the world, $\langle x_1, \ldots, x_n \rangle$.
- In Gibbs sampling, we start from a given state $s = \langle x_1, \dots, x_n \rangle$. Based on this, we generate a new state, $s' = \langle x'_1, \dots, x'_n \rangle$.
- s' depends only on s!
- ullet There is a well-defined probability of going from s to s'.

Gibbs sampling constructs a Markov chain over the Bayes net

Example



- 1. Generate a first sample: C=0, R=0, S=0, W=1.
- 2. Pick R, sample it from p(R|C=0,W=1,S=0). Suppose we get R=1.
- 3. Our new sample is C = 0, R = 1, S = 0, W = 1
- 4.

January 29, 2007

6

COMP-526 Lecture 12

Markov chains

A Markov chain is defined by:

- ullet A set of states S
- A starting distribution over the set of states $p_0(s) = p(s_0 = s)$. If the state space is discrete, this can be represented as a column vector $\mathbf{p_0}$
- ullet A stationary transition probability $p_{ss'}=p(s_{t+1}=s'|s_t=s).$ For convenience, we often put these in a n imes n matrix T

$$s_0 \rightarrow s_1 \rightarrow \cdots \rightarrow s_t \rightarrow s_{t+1} \rightarrow \cdots$$

8

Steady-state (stationary) distribution

• Where will the chain be in 1 step?

$$\mathbf{p}_{1}' = \mathbf{p}_{0}'T \longrightarrow \mathbf{p}_{1} = T'\mathbf{p}_{0}$$

where T^\prime denotes the transpose of T

• In two steps?

$$\mathbf{p_2} = T'\mathbf{p_1} = (T')^2\mathbf{p_0}$$

• In t steps?

$$\mathbf{p_t} = T' \mathbf{p_{t-1}} = (T')^t \mathbf{p_0}$$

A <u>stationary distribution</u> π is a distribution left invariant by the chain: $\pi = T'\pi$

January 29, 2007

9

COMP-526 Lecture 12

Properties of Markov chains

- Do all Markov chains converge to a stationary distribution?
 No! Consider periodic Markov chains (which contain cycles)
- When this distribution exists, is it always unique?
 No! It may depend on the initial distribution. Such chains are called reducible
- Are there conditions under which we can guarantee that the distribution is unique?
 Yes.

Properties of Markov chains

- Do all Markov chains converge to a stationary distribution?
- When this distribution exists, is it always unique?
- Are there conditions under which we can guarantee that the distribution is unique?

January 29, 2007 10 COMP-526 Lecture 12

Ergodicity

- An <u>ergodic</u> Markov chain is one in which any state is reachable from any other state, and there are no strictly periodic cycles
- In such a chain, there is a unique stationary distribution π , which can be obtained as:

$$\pi = \lim_{t \to \infty} \mathbf{p}_t$$

This is called **equilibrium** distribution

 Note that the chain reaches the equilibrium distribution regardless of p₀

 January 29, 2007
 11
 COMP-526 Lecture 12
 January 29, 2007
 12
 COMP-526 Lecture 12

Detailed balance

• Consider the stationary distribution:

$$\pi(s') = \sum_{s} \pi(s) p(s, s')$$

This can be viewed as a "flow" property: the flow out of s^\prime has to be equal to the flow coming into s^\prime from all states

 One way to ensure this is to make flow equal between <u>any pair</u> of states:

$$\pi(s)p(s,s') = \pi(s')p(s',s)$$

This gives us a <u>sufficient condition</u> for stationarity, called

detailed balance

• A Markov chain with this property is called reversible

January 29, 2007

13

COMP-526 Lecture 12

Sampling the equilibrium distribution

- We can sample π just by running the chain a long time:
 - Set $s_0 = i$ for some arbitrary i
 - For t = 1, ..., M, if $s_t = s$, sample a value s' for s_{t+1} based on p(s, s')
 - Return s_M .

If M is large enough, this will be a sample from π

 In practice, we would like to have a <u>rapidly mixing</u> chain, i.e. one that reaches the equilibrium quickly

Markov Chain Monte Carlo methods

- Suppose you want to generate samples from some distribution (but it is hard to get samples directly
 E.g., We want to sample uniformly the space of graphs with certain properties
- You set up a Markov chain such that its stationary distribution is the desired distribution
- Note that the "states" of this chain can be fairly complicated!
- You start at some state, let time pass, and then take samples
- For this to work we need to ensure:
 - that the chain has a unique stationary distribution
 - that the stationary distribution is what we want
 - that we reach the stationary distribution quickly

January 29, 2007

14

COMP-526 Lecture 12

Implementation issues

- The initial samples are influenced by the starting distribution, so they need to be thrown away. This is called the **burn-in stage**
- Because burn-in can take a while, we would like to draw several samples from the same chain!
- ullet However, if we take samples $t,\,t+1,\,t+2...$, they will be highly correlated
- Usually we wait for burn-in, then take every nth sample, for some n sufficiently large. This will ensure that the samples are (for all practical purposes) uncorrelated

16

Gibbs sampling as MCMC

- We have a set of random variables $X = \{x_1 \dots x_n\}$, with evidence variables E = e. We want to sample from p(X E|E = e).
- Let X_i be the variable to be sampled, currently set to x_i , and \bar{x}_i be the values for all other variables in $X E \{X_i\}$
- ullet The transition probability for the chain is: $p(s,s')=p(x_i'|\bar{x}_i,e)$
- Obviously the chain is ergodic
- ullet We want to show that p(X-E|e) is the stationary distribution.

Gibbs satisfies detailed balance

$$\begin{array}{lcl} \pi(s)p(s,s') & = & p(X-E|e)p(x_i'|\bar{x}_i,e) \\ \\ & = & p(x_i,\bar{x}_i|e)p(x_i'|\bar{x}_i,e) \\ \\ & = & p(x_i|\bar{x}_i,e)p(\bar{x}_i|e)p(x_i'|\bar{x}_i,e) \text{ (by chain rule)} \\ \\ & = & p(x_i|\bar{x}_i,e)p(x_i',\bar{x}_i|e) \text{ (backwards chain rule)} \\ \\ & = & p(s',s)\pi(s') \end{array}$$

January 29, 2007 18 COMP-526 Lecture 12