

Lecture10b: Gibbs Sampling

CS540 3/29/18

Announcements

On-campus students:

- Make sure I am wearing the microphone
- Check the clip!

Reading Assignment (for Tuesday)

Pedro Felzenszwalb & Daniel Huttenlocher, *Pictorial Structures for Object Recognition*, International Journal of Computer Vision 61(1): 55-79, 2005

Project #2

In project #1, most groups either

- Compared search strategies to find simple paths
- Created a two level system
 - Subgoals
 - Simple paths

For project #2

- Same simulator
- Same format for specifying problems
- But I specify the problems
- You have two weeks to improve your code
 - Anticipate hard problems
 - Parse problems into subgoals
 - Plan at level of subgoals
 - Plan paths to satisfy each subgoal

Project #2 (cont)

Two weeks from today (April 10)

- Freeze your code
- I will release my state and goal files

You have one more week (until April 17) to:

Test your system on my problems/goals

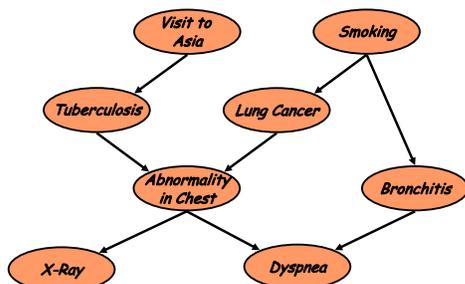
Write a 4 page paper (max) describing the results

This is a challenge paper. No motivation or literature section is required. However, it does need an introduction that summarizes what a reader should learn, a results section, an analysis of the results, and a conclusion.

Prepare a 5 minute in-class presentation (April 17)

Every team member writes one separate column describing their contribution

How do we do inference on a complex Bayesian network?



Inference by sampling

Suppose we can sample instances $\langle X_1, \dots, X_n \rangle$ according to $P(X_1, \dots, X_n)$

Want to compute $P(e)$

The probability that a random sample $\langle X_1, \dots, X_n \rangle$ satisfies e is approximately $P(e)$

We can view each sample as tossing a biased coin with probability $P(e)$ of "Heads"

BN sampling

Let X_1, \dots, X_n be order of variables consistent with arc direction

for $i = 1, \dots, n$ do
 sample x_i from $P(X_i | Pa(X_i))$
 (Note: since $Pa(X_i) \subseteq \{X_1, \dots, X_{i-1}\}$, we already assigned values to them)

return x_1, \dots, x_n

BN sampling

$P(e)$	0.001
--------	-------

$P(b)$	0.03
--------	------

$P(r e)$	0.3
$P(r \bar{e})$	0.001

$P(a b, e)$	0.98
$P(a b, \bar{e})$	0.7
$P(a \bar{b}, e)$	0.4
$P(a \bar{b}, \bar{e})$	0.01

$P(c a)$	0.8
$P(c \bar{a})$	0.05

Samples:

	B	E	A	C	R
	b	e	a	c	r

BN sampling

Sampling a complete instance is linear in number of variables

- Regardless of structure of the network

However, if $P(e)$ is small, we need many samples to get a decent estimate

Can we sample from $P(X_1, \dots, X_n | e)$?

If evidence is in roots of network, easily

If evidence is in leaves of network, we have a problem

- Our sampling method proceeds according to order of nodes in graph

Rejection sampling: keep those instantiations that are consistent with the values of the evidence variables

Estimate $P(X|e)$ by $N(X, e) / N(e)$ where $N(\cdot)$ counts the number of times an event was sampled.

Likelihood Weighting

Can we ensure that all of our samples satisfy e ?

One simple solution:

- When we need to sample a variable that is assigned value by e , use the specified value

For example: we know $Y = 1$

- Sample X from $P(X)$
- Then take $Y = 1$

Is this a sample from $P(X, Y | Y = 1)$?

Likelihood Weighting

Problem: these samples of X are from $P(X)$

Solution:

- Penalize samples in which $P(Y=1|X)$ is small

We now sample as follows:

- Let $x[i]$ be a sample from $P(X)$
- Let $w[i]$ be $P(Y = 1 | X = x[i])$

$$P(X = x | Y = 1) \approx \frac{\sum_i w[i] P(X = x | x[i])}{\sum_i w[i]}$$

Markov chain Monte Carlo sampling

Generates events by making random changes to the state variable.

The next state is generated by sampling a value for one of the nonevidence variables conditioned on the current values.

Markov chains

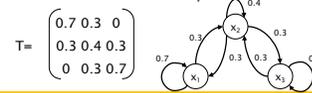
A Markov chain is a random process (infinite sequence of random variables)

$(X(0), X(1), \dots, X(t), \dots)$ that satisfies:

$$P(X(t) | X(0), \dots, X(t-1)) = P(X(t) | X(t-1))$$

The probability of a particular state at time t depends only on the state at time $t-1$

If the transition probabilities are fixed for all t , the chain is called *homogeneous* and is characterized by a transition matrix T .



Sampling via Markov chains

For sampling from $P(x)$, we require that for any starting state $x(0)$:

$$\lim_{t \rightarrow \infty} P_t(x) = P(x)$$

Equivalently, the stationary distribution of the Markov chain must be $P(x)$:

$$P_{t+1}(x') = \sum_x P_t(x) Q(x \rightarrow x')$$

the transition probability from x to x'

$$P(x') = \sum_x P(x) Q(x \rightarrow x')$$

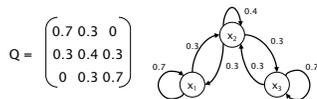
Using a Markov chain to sample

If a Markov chain converges to is probability distribution:

- Start from an arbitrary state
- Do a long random walk
- According to the transition probabilities
- Count how often you are in each state x
- Divide the occurrences of x by t to get $p(x)$

The resulting state will be sampled from $P(x)$.

Stationary distribution



The stationary distribution of this chain is $(0.33, 0.33, 0.33)$

Markov chains for sampling

To ensure that the chain converges to a unique stationary distribution the following conditions are sufficient:

- *Irreducibility*: every state is eventually reachable from any start state; for all x, y there exists a t such that $P_t(y) > 0$ when starting at x
- *Aperiodicity*: the chain doesn't get caught in cycles.

The process is *ergodic* if it is both irreducible and aperiodic

Detailed balance

To ensure that the stationary distribution of the Markov chain is $P(x)$ it is sufficient for P and Q to satisfy the *detailed balance (reversibility)* condition:

$$P(x)Q(x \rightarrow x') = P(x')Q(x' \rightarrow x)$$

Given that detailed balance holds:

$$\begin{aligned} \sum_x P(x)Q(x \rightarrow x') &= \sum_x P(x')Q(x' \rightarrow x) \\ &= P(x') \sum_x Q(x' \rightarrow x) \\ &= P(x') \end{aligned}$$

Back to Bayesian Networks

What does this have to do with inference on a Bayesian Network?

- Bayesian nets have random variables (nodes)
- Bayesian nets have relations (edges)
- Think of a complete assignment of values to random variables at a time t as a state $x(t)$ in our Markov chain.
- The relations tell us how to probabilistically move from one state to another

Gibbs sampling

Idea: To transition from one state (variable assignment) to another by:

- Pick a variable X_j ,
- Sample its value from the conditional distribution

$$P(x_j \mid x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_n)$$

In a Bayesian network x_j depends only on a subset of the variables.

Markov Blanket

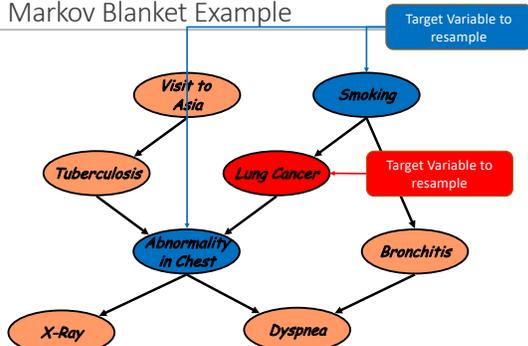
Variables are independent of their non-descendants given their parents

Variables are independent of *everything else in the network* given their *Markov blanket*.

So, to sample a node, only need to condition on its Markov blanket:

$$\begin{aligned} P(x_j \mid x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_n) &= \\ P(x_j \mid MB(x_j)) \end{aligned}$$

Markov Blanket Example



The Gibbs sampling algorithm

$GIBBS(X, e, bn, N)$ returns estimate of $P(X|e)$

$N[x]$ - counts the number of times each value of X was observed

$x[j]$ - the current state of the network $x[0]$ initialized with random values for the nonevidence variables

for $j = 1$ to N do

for each nonevidence variable X_i

sample X_i from $P(X_i|MB(X_i))$

$N[x] = N[x] + 1$, where x is the value of X in $x[j]$

Convergence of Gibbs sampling

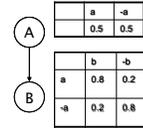
Gibbs sampling satisfies detailed balance:

$$\begin{aligned}
 P(\mathbf{x}|\mathbf{e})P(x'_i|\bar{\mathbf{x}}_i, \mathbf{e}) &= P(x_i, \bar{\mathbf{x}}_i|\mathbf{e})P(x'_i|\bar{\mathbf{x}}_i, \mathbf{e}) \\
 &= P(x_i|\bar{\mathbf{x}}_i, \mathbf{e})P(\bar{\mathbf{x}}_i, \mathbf{e})P(x'_i|\bar{\mathbf{x}}_i, \mathbf{e}) \\
 &= P(x_i|\bar{\mathbf{x}}_i, \mathbf{e})P(x'_i, \bar{\mathbf{x}}_i|\mathbf{e}) \\
 &= P(\mathbf{x}'|\mathbf{e})P(x_i|\bar{\mathbf{x}}'_i, \mathbf{e})
 \end{aligned}$$

$\bar{\mathbf{x}}_i$

Gibbs sampling example

Consider a 2 variable network:



Initialize randomly

Sample variables alternately

Practical issues

How many iterations?

When to stop?