

① Approximate Inference w/ Sampling.

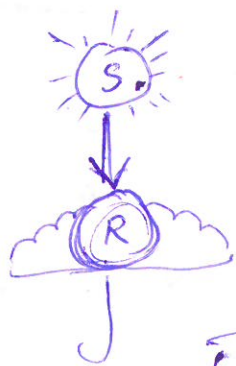
(i.1) Motivation: Exact inference is intractable for many types of BNs.

(1.2) Idea:

- > A BN is a generative model.
- > You can generate samples from it.
- > Use the samples to estimate the probabilities you want.

(1.3) Simple Example:

$$\sum_R P(R|S) = 1$$



| S | P(S) |
|---|------|
| | 0.8 |
| | 0.2 |

| S | R | P(R S) |
|---|---|--------|
| | | 0.85 |
| | | 0.15 |
| | | 0.20 |
| | | 0.80 |

> w/ these coins you can generate samples:

| | P() |
|----|---------------------------------|
| 1. | $\frac{\#(\text{Umbrella})}{7}$ |
| 2. | |
| 3. | $P(\text{Umbrella})$ |
| 4. | |
| 5. | $1 - P(\text{Umbrella})$ |
| 6. | |
| 7. | |

> Query: ~~$P(R|\text{Sun})$~~

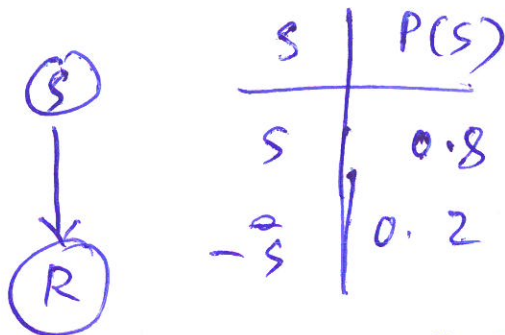
$P(R) = ?$ ~~$P(R=\text{Umbrella})$~~ or ~~$P(R=\text{No Umbrella})$~~

> $P(\text{Sun}) = 0.8$

$P(\text{Cloud}) = 0.2$

Flip this coin: Then use sunny coin.
else use cloudy coin.

> $P(\text{Umbrella}) = 0.85$ & $P(\text{Umbrella}) = 0.20$



| | S | R | P(R S) |
|---|----|----|--------|
| ① | S | x | 0.85 |
| | S | -x | 0.15 |
| ② | -S | x | 0.2 |
| | -S | -x | 0.8 |

> Flip a coin for S

S

> Flip the coin ① to obtain value for R

-x

(S, -x)

(-S, -x)

(S, x)

⋮

⋮

(S, -x)

use coin ② to get R value.

$P_x(S, x)$

$P_x(S=S)$

$P_x(R=x)$

2.4) Are these probability estimates "correct"?

Consistency: A sampling estimate is consistent if in the limit it yields the exact probability.

i.e.) if $\hat{P}(x | E=e) \approx \#_s(x, e) / \#_s(e)$ is the sampling estimate

Then \hat{P} is consistent if

$$\lim_{N \rightarrow \infty} \hat{P}(x | E=e) = P(x | E=e)$$

$$\text{B. i.e.)} = \lim_{N \rightarrow \infty} \frac{\#_s(x, e)}{N} \cdot \frac{1}{\frac{\#_s(e)}{N}}$$

Straightforward

Because of how we defined the sampling process

$$\begin{aligned} \lim_{N \rightarrow \infty} \frac{\#_s(x, e)}{N} &= P(x, e) \\ \lim_{N \rightarrow \infty} \frac{\#_s(e)}{N} &= P(e) \end{aligned}$$

$$\begin{aligned} \Rightarrow &= \frac{P(x, e)}{P(e)} \\ &= P(x | e) \end{aligned}$$

Introduction

3.1 Algorithm : Input: $Q, E=e, H = [X = Q, H]$

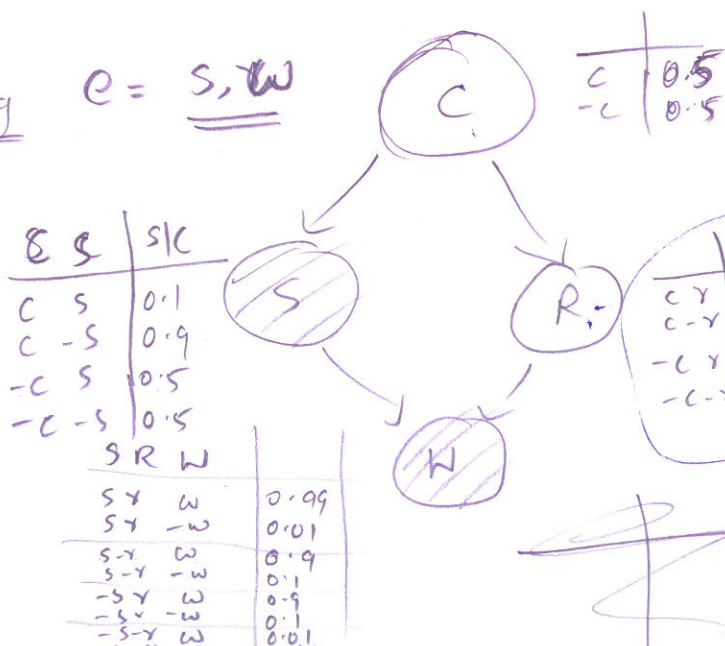
1. Fix $E=e$
2. Topologically sort Q, E, H
3. In sorted order, sample values for each var: $P_r(X_i | \text{Parents}(X_i))$
 ↳ evidence vars are fixed to i/ps.

4. For each sample $(q, e, h) = s_j$ weight it by $\prod P_r(e_i | \text{Parents}(E_i))$

$Wt_j \leftarrow \prod_{E_i \in e} P_r(E_i | \text{Parents}(E_i))$

5. $P_r(q | e) = \frac{\sum_{s_j \in S_q} Wt_j}{\sum_{s_j \in S} Wt_j}$

e.g $e = s, w$



| Samples | Weight |
|-------------------------------|-------------------|
| $P_r(S C) \times P_r(W S, C)$ | |
| $C \ S \ Y \ W$ | 0.1×0.99 |
| $-C \ S \ -Y \ W$ | 0.5×0.9 |
| $C \ S \ -Y \ W$ | |
| $-C \ S \ Y \ W$ | |
| $C \ -S \ Y \ W$ | |
| $-C \ -S \ -Y \ W$ | |
| $C \ -S \ -Y \ W$ | |
| $-C \ -S \ Y \ W$ | |
| $C \ -S \ Y \ -W$ | |
| $-C \ -S \ -Y \ -W$ | |
| $C \ -S \ -Y \ -W$ | |
| $-C \ -S \ Y \ -W$ | |
| $C \ S \ Y \ -W$ | |
| $-C \ S \ Y \ -W$ | |
| $C \ S \ -Y \ -W$ | |
| $-C \ S \ -Y \ -W$ | |
| $C \ -S \ Y \ -W$ | |
| $-C \ -S \ Y \ -W$ | |
| $C \ -S \ -Y \ -W$ | |
| $-C \ -S \ -Y \ -W$ | |

$P_r(C | S, W) = \frac{\sum Wt_j}{\text{Total wts.}}$

4

Monte Carlo Markov Chain

Motivation:

Idea: Let's ensure that evidence influences both upstream and downstream.

Idea: Starting from ^{top} Scratch is the Problem

~~Why don't we keep~~

- > Sample ~~every~~ variables conditioned on rest.
- > No re-start from scratch/top.

Gibbs Sampling

- > start w/ an arbitrary instantiation but consistent w/ evidence.

x_1, x_2, \dots, x_n
" " " "
 x_1, x_2, \dots, x_n

- > select one variable (keep evidence) x_i
fixed

- > Sample a ^{new} value conditioned on the $P(x_i | x_{-/i})$
rest

↑
all but x_i

- > Repeat!

8

BN Summary

- > BN compactly specifies joint distributions in terms of CPTs.
- > Any prob. query can be evaluated over a BN!
 - ↳ Remarkable since we only store 'n' CPTs.
- > Inference in BNs is the key challenge
 - >> Exact inference is hard
 - >>> For some n/ws is tractable.
 - >> Sampling, approx inference, is useful in practice.
- > How to obtain the structure of n/w?
 - ↳ domain experts
 - ↳ some structure learning
- > How to obtain CPTs?
 - ↳ hard-code \rightarrow not tenable
 - ↳ learn (see example in following lect).

Oct 24th Lecture

Plan

Inference
in BN

- (1) Recap of Sampling based inference
- (2) MCMC : Gibbs sampling
- (3) Bayesian Network summary

Learning

- (1) Learning from examples
- (2) Basic supervised learning recipe
- (3) Naïve Bayes
- (4) Logistic Regression

$$Q = Q_1, Q_2, Q_3$$

$$\underline{Q} = \begin{matrix} Q_1, -Q_2, \\ Q_1, Q_2, -Q_3 \end{matrix}$$

$$\underline{Q_1, -Q_2, -Q_3}$$

$$\underline{Q_1, -Q_2, Q_3}$$

$$\underline{Q_1, -Q_2, Q_3}$$

①

Sampling Based Inference - Recap

①

> Exact inference is intractable for general BNs

$O \rightarrow O \rightarrow O \dots$ > VE is tractable - polynomial $O(d^2 n)$ for chains. More generally for poly-tree
size of the dom
the # of vars

Motivation

> When \exists more than one path to any node, inference is intractable - exponential in BN size.

> VO helps but isn't enough.

②

> obtain samples from BN.

> Use samples to estimate probabilities

Idea

③

> Suppose there are samples:

| B | A | C | E |
|----|----|----|----|
| b | a | -c | e |
| -b | a | c | -e |
| b | -a | c | -e |

What is $Pr(c)$?

What is $Pr(a|-c)$?

$$Pr(c) = \frac{\#(c)}{N} = \frac{2}{3} \quad Pr(-c) = \frac{1}{3}$$

$$Pr(a|-c) = \frac{\#(a, -c)}{\#(-c)} = \frac{1}{1} = 1$$

④

How do we obtain samples?

Question

⑤ - cont'd

> Likelihood weighted sampling



















↳ Don't sample evidence vars

↳ keep their values fixed

↳ sample the rest

Beware: Probs are inconsistent

> If E is downstream from Q

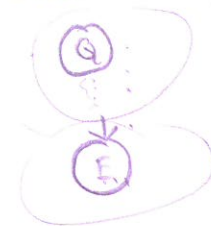
| Count | E | Q |
|--------------------------|---|---|
| $w(s_1) = Pr(e q_1) s_1$ |  |  |
| $w(s_2) = Pr(e q_2) s_2$ |  |  |
| s_3 |  |  |
| |  |  |
| |  |  |
| |  |  |
| |  |  |
| |  |  |
| $w(s_n) = Pr(e q_n) s_n$ |  |  |

> Q samples don't account for evidence E

↳ ok if E is the top of the n/w



> But not necessarily true in general



> Idea?

Formally

$$> \underline{Pr(q_i | e)} = \frac{\sum_{s_i \in S_{q_i}} \underline{w(s_i)}}{\sum_{s_i \in S_{q_i}} w(s_i) + \sum_{s_i \in S_{-q_i}} w(s_i)}$$

$$> \underline{w(s_i)} = \underline{Pr(e | q_i)}$$

$$\sum_{s_i \in S_{q_i}} w(s_i) + \sum_{s_i \in S_{-q_i}} w(s_i)$$

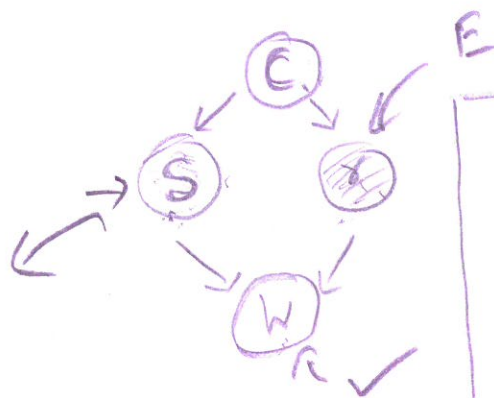
⑥ Markov Chain Monte Carlo (MCMC)

Motivation

Evidence should influence both
up & down stream vars!

The issue

S sampled
from $P(S|C)$



> we want S to be
influenced by ~~X~~

> Likelihood weighted doesn't
directly do this

Why?

Starting to sample from
top every time decouples S from X

The fix

Not from
the top \Rightarrow

When sampling value for
a var condition on the
current values for all
vars in the n/w

Why the fix works?

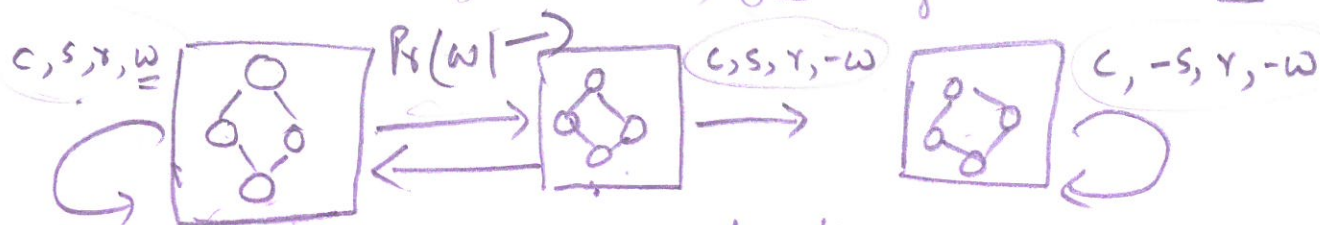
Suppose you have sampled
C=c, W=w & R=r \rightarrow Evidence

If you sample S from $P(S|C=c, W=w, R=r)$
you will see that S influenced by r

⑦ - contd

Why does Gibbs work?

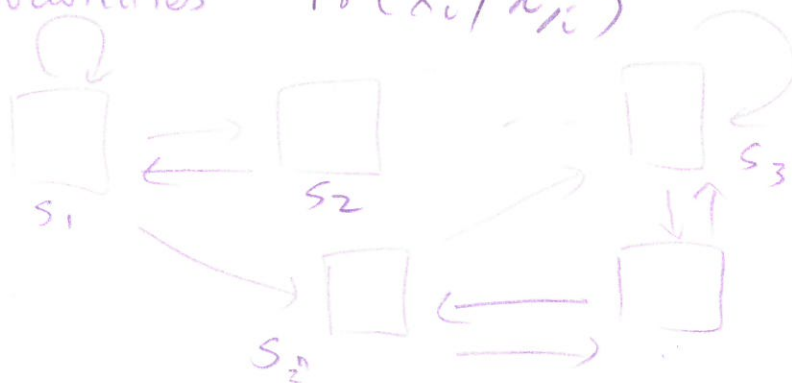
> Gibbs cycles through a Markov chain over diff configs of the vars



Markov ~~chain~~ ^{network} of config.

> Each configuration = sample = state in MC

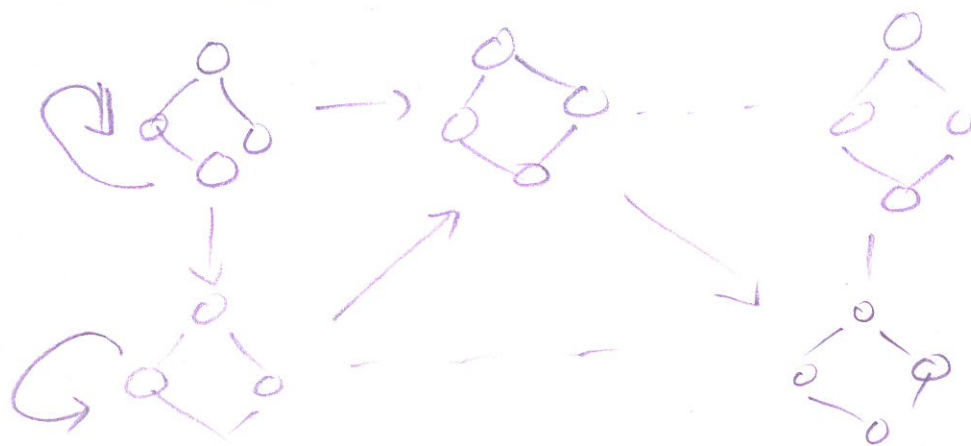
> A random walk on this chain cond. on the probabilities $Pr(x_i/x_{-i})$



> If you walk forever then stationary distribution $\pi \rightarrow Pr(s_i)$, the prob. of sample according to the BN

Why does Gibbs work?

- > The set of configurations that Gibbs cycles through forms a Markov chain.



$2^4 = 16$
states

- > Each state = sample.
- > Walk through this chain = seq of samples.
- > If you walk forever, the walk is transition guided by the $Pr(x_i | x_{-i})$ distribution.
- > One can show that the stationary distribution of this walk (ie) the prob of landing in each state = $Pr(\text{sample})$ given the BN!
- > so we get consistent estimates.