

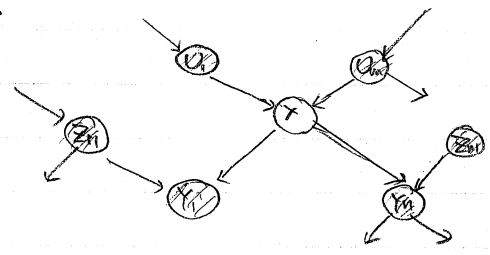
10/12 Approximate inference

Last class : 1) Rejection sampling (slow)  
2) Likelihood weighting (faster)

Today : 3) Markov chain Monte Carlo (MCMC) (fastest)

Def : Markov blanket  $B_x$  of node  $X$  consists of parents, children, and spouses of  $X$ .

Thm :  $X$  is conditionally independent of all nodes outside  $B_x$  given nodes in  $B_x$ .

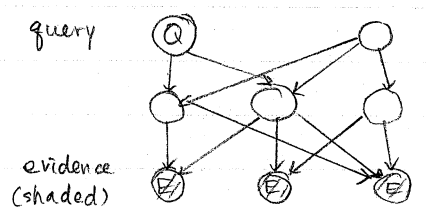


MCMC stimulation

Query nodes  $Q$ .

Evidence nodes  $E$ .

How to estimate  $P(Q=q | E=e)$ ?



\* To estimate  $P(Q=q | E=e)$ :

- fix evidence nodes to observed values.
- initialize non-evidence nodes at random.
- repeat  $N$  times:
  - pick non-evidence node  $X \notin E$  at random.
  - use Bayes rule to compute  $P(X | B_x)$  where  $B_x$  is fixed to current values.
  - resample  $X$  from  $P(X=x | B_x)$

\* Count # samples  $N(q)$  where  $Q=q$

\* Estimate  $P(Q=q | E=e) \approx N(q)/N$ .

Converges in limit  $N \rightarrow \infty$  to correct answer.

\* Key difference between likelihood weighting (LW) and MCMC:

LW > sample non-evidence nodes from  $P(X | pa(X))$   
 MCMC > sample non-evidence nodes from  $P(X | B_x)$

**Learning**

\* BN = DAG + CPTs not always available from experts.

How to learn from examples?

\* Maximum likelihood (ML) estimation.

- simplest form of learning in BNs.

- choose ("estimate") model (DAG + CPTs) to maximize

$P(\text{observed data} \mid \text{DAG} + \text{CPTs})$   
 likelihood.

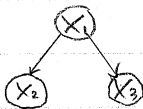
Case I known structure of DAG, lookup tables for CPTs, complete data

- DAG is fixed over some known, finite set of discrete nodes  $\{X_1, X_2, \dots, X_n\}$

- CPTs enumerate  $P(X_i = x \mid pa(X_i) = \pi)$  as lookup table  
 ← parent configuration.

- Data is T complete instantiations of nodes in BN.

Ex:



example #	$X_1$	$X_2$	$X_3$
1	1	0	1
2	0	0	0
3	0	1	0
⋮			
T	1	0	1

jargon: "complete data", "fully observed", "no hidden nodes", "fully visible".

More generally, denote data as  $\prod_{t=1}^T (X_1^{(t)}, X_2^{(t)}, \dots, X_n^{(t)})^T$ .

\* IID assumption

Samples are independently identically distributed from joint distribution

$P(X_1, X_2, \dots, X_n)$  of BN.

\* Probability of IID data set.

$$P(\text{data}) = \prod_{t=1}^T P(X_1 = x_1^{(t)}, X_2 = x_2^{(t)}, \dots, X_n = x_n^{(t)}) \text{ due to IID assumption.}$$

(product over rows)

Probability of t-th example:

$$P(X_1 = x_1^{(t)}, X_2 = x_2^{(t)}, \dots, X_n = x_n^{(t)}) = \prod_{i=1}^n P(X_i = x_i^{(t)} \mid X_1 = x_1^{(t)}, X_2 = x_2^{(t)}, \dots, X_{i-1} = x_{i-1}^{(t)})$$

← product rule

$$= \prod_{i=1}^n P(X_i = x_i^{(t)} \mid pa(X_i) = pa_i^{(t)})$$

← conditional independence from DAG.

\* log-likelihood  $\mathcal{L}$

$$\mathcal{L} = \log P(\text{DATA})$$

$$= \log \prod_{t=1}^T P(X_1^{(t)}, X_2^{(t)}, \dots, X_n^{(t)}) \quad \text{IID}$$

$$= \log \prod_{t=1}^T \prod_{i=1}^n P(X_i^{(t)} | pa_i^{(t)}) \quad \text{due to prod rule \& CI.}$$

$$= \sum_{t=1}^T \sum_{i=1}^n \log P(X_i^{(t)} | pa_i^{(t)})$$

$$= \sum_{i=1}^n \sum_{t=1}^T \log P(X_i^{(t)} | pa_i^{(t)}) \quad \text{swapping order of summation.}$$

Let  $\text{count}(X_i = x, pa_i = \pi)$  denote # examples for which  $X_i = x$  and  $pa_i = \pi$ .

$$\mathcal{L} = \sum_{i=1}^n \sum_x \sum_{\pi} \underbrace{\text{count}(X_i = x, pa_i = \pi)}_{\substack{\text{possible} \\ \text{values of } X_i}} \log \underbrace{P(X_i = x | pa_i = \pi)}_{\substack{\text{properties of data} \\ \text{unknowns to be optimized} \\ \text{(learned from data)}}$$

Write  $\mathcal{L} = \sum_{i, \pi} \mathcal{L}_{i\pi}$  where  $\mathcal{L}_{i\pi} = \sum_x \text{count}(X_i = x, pa_i = \pi) \log P(X_i = x | pa_i = \pi)$

From decomposition  $\mathcal{L} = \sum_{i, \pi} \mathcal{L}_{i\pi}$ , we can independently optimize CPT entries at each node in BN and for each parent configuration of that node.

\* ML estimation

For each node  $X_i$ , for each row  $\pi$  of CPT, maximize  $\mathcal{L}_{i\pi}$  subject to:

$$(i) \sum_x P(X_i = x | pa_i = \pi) = 1$$

$$(ii) P(X_i = x | pa_i = \pi) \geq 0$$

\* Shorthand notation. at node  $i$ , row  $\pi$  of CPT:

$$\text{let } C_x = \text{count}(X_i = x, pa_i = \pi)$$

$$\text{let } p_x = P(X_i = x | pa_i = \pi)$$

$$\text{how to maximize } \sum_x C_x \log p_x \quad \text{subject to } \begin{cases} \sum p_x = 1 \\ p_x \geq 0 \end{cases} ?$$

$$\text{Solution: } p_x = \frac{C_x}{\sum_{\beta} C_{\beta}}$$

$$\begin{aligned} \text{ML Solution: } P_{\text{ML}}(X_i = x | pa_i = \pi) &= \frac{\text{count}(X_i = x, pa_i = \pi)}{\sum_{x'} \text{count}(X_i = x', pa_i = \pi)} \\ &= \text{count}(X_i = x, pa_i = \pi) / \text{count}(pa_i = \pi) \end{aligned}$$

\* Properties

- Asymptotically correct

$$P_{\text{ML}}(X_1, X_2, \dots, X_n) \rightarrow P(X_1, X_2, \dots, X_n) \text{ as } T \rightarrow \infty$$

- Problematic in non-asymptotic regime. (of "sparse" data)

$$P_{\text{ML}}(X_i = x | pa_i = \pi) = \begin{cases} 0 & \text{if } \text{count}(X_i = x, pa_i = \pi) = 0 \\ \text{undefined} & \text{if } \text{count}(pa_i = \pi) = 0 \end{cases}$$

Ex: Markov models of language

\* let  $w_l$  denote  $l^{\text{th}}$  word in sentence (or utterance or paragraph)

How to model  $P(w_1, w_2, \dots, w_L)$ ?

\* Simplifying assumptions

(i)  $P(w_l | w_1, w_2, \dots, w_{l-1}) = P(w_l | w_{l-c_{l-1}}, w_{l-c_{l-2}}, \dots, w_{l-1})$  finite context.

(ii)  $P(w_l = w | w_{l-c_{l-1}} = v_{l-1}, \dots, w_{l-1} = v_1) = P(w_{l+1} = w | w_{l+1-c_{l-1}} = v_{l-1}, \dots, w_l = v_1)$   
position invariance