

Review

\* Inference in BNs

evidence nodes  $E$

query node  $Q$

How to compute  $P(Q|E)$ ?

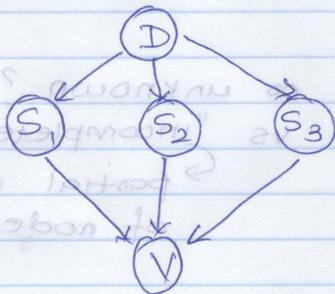
\* Polytrees

- singly connected networks

- polynomial time inference.

\* Loopy BNs

One approach: node clustering



disease  $D \in \{0, 1\}$

visit  $V \in \{0, 1\}$



Merge nodes to form polytree



Merge  $S_1, S_2, S_3$  into "mega-node"



Merge CPTs into "mega-CPTs".

take on  $2^3$  values

$S_1$	$S_2$	$S_3$	$S$	$P(S D)$
0	0	0	0	$P(S_1=0 D) P(S_2=0 D) P(S_3=0 D)$
0	0	1	1	$P(S_1=0 D) P(S_2=0 D) P(S_3=1 D)$
...	...	...	...	...
1	1	1	1	$P(S_1=1 D) P(S_2=1 D) P(S_3=1 D)$

\* Polynomial Polytree algorithm linear in size of CPTs, but CPTs growing exponentially with clustering.

\* How to choose optimal clustering? computationally hard problem.

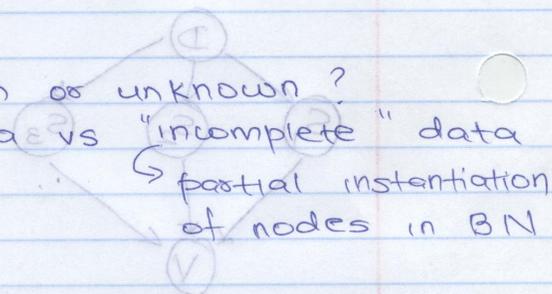
**Learning**

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

How to learn from examples?

\* Issues

- structure (DAG) - known or unknown?
- evidence: complete data vs "incomplete" data



- optimization:
  - combinatorial vs ~~iterative~~ continuous
  - (eg: learning DAG) (eg: learning CPTs)
- algorithms: non-iterative vs iterative (loop many times over data set)
- solution: local vs global optima in model estimation.

\* Maximum likelihood estimation (ML)

- simplest form of learning  
 (choose ("estimate") model (DAG + CPTs) to

maximize  $P(\text{observed data} | \text{model})$   
 "likelihood"

$P(1=1|D) P(2=2|D) P(3=2|D) P(4=1|D)$

**Ex : biased coin**

$X \in \{heads, tails\}$

$P(X=heads) = p$   
 $P(X=tails) = 1-p$

\* How to estimate  $p$  from  $T$  examples (i.e., results of  $T$  coin tosses)?

\* I.I.D. assumption

Samples are independently, identically distributed according to  $P(X)$ .

$\rightarrow \{X^{(1)}, X^{(2)}, \dots, X^{(T)}\}$   $T$  samples

\* Probability of I.I.D. data:

$P(\text{data}) = P(X = x^{(1)}, X = x^{(2)}, X = x^{(3)}, \dots, X = x^{(T)})$   
 $= P(X = x^{(1)}) P(X = x^{(2)}) \dots P(X = x^{(T)})$  (i.c. coin tosses independent)  
 $= \prod_{t=1}^T P(X = x^{(t)})$

\* Log-probabilities

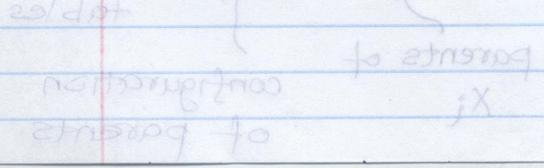
$\mathcal{L} = \log P(\text{data})$

log-likelihood

$= \log \prod_{t=1}^T P(X = x^{(t)})$

$= \sum_{t=1}^T \log P(X = x^{(t)})$

\* CPTs: enumerate  $P(x_i = x | \text{parents}(x_i))$



\* Notation :

Let  $N_H = \text{count}(X = \text{heads})$

$N_T = \text{count}(X = \text{tails})$

Clearly  $N_H + N_T = T$  (total # samples)

In terms of counts :

$$\mathcal{L} = N_H \log p + N_T \log(1-p)$$

\* Maximum likelihood estimation :

$$\frac{\partial \mathcal{L}}{\partial p} = \frac{N_H}{p} + \frac{N_T}{1-p} (-1) = 0 \quad \text{at maximum}$$

$$N_H(1-p) - N_T p = 0$$

$$N_H - p(N_H + N_T) = 0$$

$$p = \frac{N_H}{N_H + N_T} = \frac{N_H}{T}$$

- intuitively, ML estimate is relative empirical frequency of heads.

Discrete BNs with "complete data"

\* Given : fixed DAG over discrete nodes  $\{X_1, X_2, \dots, X_n\}$

\* CPTs enumerate  $P(X_i = x \mid \text{pa}(X_i) = \pi)$  as lookup tables

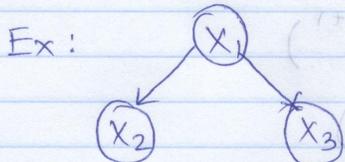
parents of  $X_i$

configuration of parents

conditional dependence in BN

\* Data is  $T = \{ \text{complete instantiations of all nodes in BN.} \}$

$$\left\{ (X_1^{(t)}, X_2^{(t)}, X_3^{(t)}, \dots, X_n^{(t)}) \right\}_{t=1}^T$$



$$X_i \in \{0, 1\}$$

$$n = 3$$

t	X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>
1	1	0	1
2	1	1	1
3	0	1	1
4	0	1	0
⋮			
T	0	1	0

\* Each n-tuple of values is called an "example"

Goal: learn from examples

estimate CPTs  $P(X_i = x_i | \text{pa}_i = \pi)$

that maximize probability of data  
likelihood

\* I.I.D. assumption

Examples are independently, identically

distributed according to  $P(X_1, X_2, \dots, X_n)$

\* Probability of data

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

← Joint prob of  $t^{\text{th}}$  example

\* Work out  $t^{\text{th}}$  term: Product Rule

$$P(X_1 = x_1^{(t)}, \dots, X_n = x_n^{(t)}) = P(X_1 = x_1^{(t)}) P(X_2 = x_2^{(t)} | X_1 = x_1^{(t)})$$

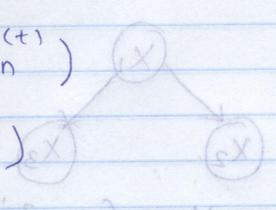
$$= \prod_{i=1}^n P(X_i = x_i^{(t)} | X_1 = x_1^{(t)}, \dots, X_{i-1} = x_{i-1}^{(t)})$$

## conditional dependence in BN

$$l(\theta) = \prod_{i=1}^n P(X_i = x_i^{(t)} | pa(X_i) = \pi_i^{(t)})$$

### \* Log-Likelihood

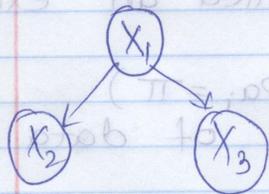
$$\begin{aligned} \mathcal{L} &= \log P(\text{data}) \\ &= \log \prod_{t=1}^T P(X_1^{(t)}, X_2^{(t)}, \dots, X_n^{(t)}) \\ &= \log \prod_{t=1}^T \prod_{i=1}^n P(X_i^{(t)} | pa_i^{(t)}) \\ &= \sum_{t=1}^T \sum_{i=1}^n \log P(X_i^{(t)} | pa_i^{(t)}) \end{aligned}$$



$$\mathcal{L} = \sum_{i=1}^n \sum_{t=1}^T \log P(X_i = x_i^{(t)} | pa(X_i) = \pi_i^{(t)})$$

swapped order of sums

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



t	X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>	count(X <sub>3</sub> =1, X <sub>1</sub> =0) = 2
1	0	0	0	count(X <sub>1</sub> =1) = 3
2	1	1	0	count(X <sub>2</sub> =1, X <sub>1</sub> =1) = 1
3	1	0	1	
4	1	0	0	
5	0	0	1	

\* CPTs that we can choose

### \* Log-likelihood

$$\mathcal{L} = \sum_{i=1}^n \sum_{x \in \mathcal{X}_i} \sum_{\pi \in \mathcal{P}_i} \text{count}(X_i = x, pa_i = \pi) \log P(X_i = x | pa_i = \pi)$$

↑ values of  $X_i$ 
↑ values of  $pa(X_i)$ 
↑ completely determined by data

\* ML estimation

How to choose  $P(X_i = x | p_{a_i} = \pi)$  to maximize  $\mathcal{L}(\text{data})$ ?

\* ML solution (w/o proof)

$$\begin{aligned} P_{ML}(X_i = x | p_{a_i} = \pi) &= \frac{\text{count}(X_i = x, p_{a_i} = \pi)}{\sum_{x'} \text{count}(X_i = x', p_{a_i} = \pi)} \\ &= \frac{\text{count}(X_i = x, p_{a_i} = \pi)}{\text{count}(p_{a_i} = \pi)} \end{aligned}$$

\* Properties of MLE

- Asymptotically correct:  $P_{ML}(X_1, X_2, \dots, X_n) \rightarrow P(X_1, X_2, \dots, X_n)$  as  $T \rightarrow \infty$

with enough data, you recover the true model

- Problematic for sparse data:

$$P_{ML}(X_i = x | p_{a_i} = \pi) = 0 \text{ if } \text{count}(X_i = x, p_{a_i} = \pi) = 0$$

$$P_{ML}(X_i = x | p_{a_i} = \pi) \text{ undefined if } \text{count}(p_{a_i} = \pi) = 0$$