# CSE 150 – 1/27/10

## Review

\* Inference in BNs
  evidence node E
    query node Q

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

\* Poly trees
  - singly connected networks
  - polynomial time inference

\* Loopy BNs

  Exact inference : node clustering
  Approximate inference: stochastic simulation
  
  [covered later    ]

## 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
                    ↳ partial instantiation of nodes in BN
  - optimization:
    combinatorial     vs.   continuous
    (e.g. learning DAG)        (e.g. learning CPTs)
  - algorithms: non-iterative   vs.  iterative
                              (loop over data many times)
  - solution: local vs. global optimum.

(1)

* Maximum Likelihood (ML) estimation
   – simplest form of learning in ~~XXXX~~ BNs
   – choose ("estimate") the model (DAG + CPTs)
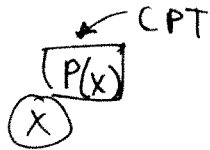     to maximize $\underbrace{P(\text{observed data}|\text{model})}$
     
     "likelihood"

Ex: biased coin

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

$P(X = \text{heads}) = P$

$P(X = \text{tails}) = 1 - P$

Trivial BN

CPT
$\boxed{P(x)}$
$(X)$

* How to estimate $p$ from observed samples
   (results of ~~XXX~~ T coin tosses)?

* IID assumptions
   Samples are $\underline{i}$dependently, $\underline{i}$dentically $\underline{d}$istribute to $P(x)$.
   $\rightarrow \{X^{(1)}, X^{(2)}, \dots, X^{(T)}\}$  T samples

* Probability of IID data:
   $$P(\text{data}) = P(X = x^{(1)}) P(X = x^{(2)}) \cdots P(X = x^{(T)})$$
   $$= \prod_{t=1}^{T} P(X = x^{(t)})$$

* Log-probability $\mathcal{L}$
   $$\mathcal{L} = \log(P(\text{data})) = \log \prod_{t=1}^{T} P(X = x^{(t)}) = \sum_{t=1}^{T} \log P(X = x^{(t)})$$
   "log-likelihood"

Let $N_H$ = count $(X = heads)$
Let $N_T$ = count $(X = tails)$

Clearly: $N_H + N_T = T \leftarrow$ total samples

In terms of counts:
$$\mathcal{L}(p) = N_H \log(p) + N_T \log(1-p)$$

\* Maximum Likelihood estimation

$$\frac{\partial \mathcal{L}}{\partial p} = \frac{N_H}{p} + \frac{N_T}{1-p} \cdot (-1) = 0$$

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

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

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

intuitively, maximum likelihood estimate of $p = P(X=heads)$ is relative frequency in observed coin ~~tosses.~~ tosses.
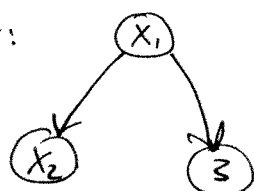
---

| Discrete BNs with "complete data" |
|---|

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

\* CPTs enumerate $P(X_i = x_i \mid pa(x_i) = \pi)$ as look up tables
$\qquad\qquad\qquad\quad\underset{\text{parents of } X_i}{\uparrow} \quad \underset{\text{parent configuration}}{\uparrow}$

\* Data is $T$ complete instantiations of nodes in BN
$$\{x_1^{(t)}, x_2^{(t)}, \dots, x_n^{(t)}\}_{t=1}^{T}$$

Ex:

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

$n = 3$

Data

| $t$th sample | $X_1$ | $X_2$ | $X$ |
|---|---|---|---|
| 1 | 0 | 0 | 0 |
| 2 | 1 | 1 | 0 |
| 3 | 1 | 1 | 0 |
| ⋮ | | | |
| $T$ | 0 | 1 | 3 |

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

   Goal: learn from examples;

      estimate CPTs $P(X_i = x \mid pa_i = \pi)$

      that maximize <u>probability</u> of data set

             likelihood

\* I.I.D. Assumption

  samples are independently, identically distributed

  according to $P(X_1, X_2, \ldots X_n)$.

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

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

           ↖ probability of $t^{th}$ example

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

$$P(X_1 = x_1^{(t)}, \ldots, X_n = x_n) = P(X_1 = x_1^{(t)}) P\left(X_2 = x_2^{(t)} \mid P(X_1 = x_1^{(t)})\right) \times \ldots \quad \text{product rule}$$

$$= \prod_{t=1}^{n} P\left(X_i = x_i^{(t)} \mid X_1 = x_1^{(t)}, \ldots, X_{i-1} = x_{i-1}^{(t)}\right)$$

$$= \prod_{t=1}^{n} P\left(X_i = x_i \mid pa(x_i) = pa_i^{(t)}\right) \quad \text{conditional dependence}$$

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

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

$$= \log \prod_{t=1}^{T} P\left(x_1^{(t)}, x_2^{(t)}, \ldots, x_n^{(t)}\right)$$

$$= \log \prod_{t=1}^{T} \prod_{i=1}^{n} P\left(x_i^{(t)} \mid pa(x_i) = pa_i^{(t)}\right)$$

$$= \sum_{t=1}^{T} \sum_{i=1}^{n} \log P\left(X_i = x_i^{(t)} \mid pa(x_i) = pa_i^{(t)}\right)$$

$$= \sum_{i=1}^{n} \sum_{t=1}^{T} \log P\left(X_i = x_i^{(t)} \mid pa(x_i) = pa_i^{(t)}\right) \quad \text{swap order of sums}$$

④

* Let $\text{count}(X_i = x, pa_i = \pi)$ denote examples for which $X_i = x_*$ and ~~pPPqPP~~ $pa(X_i) = \pi$.

Data set

| t | $X_1$ | $X_2$ | $X_3$ | $\cdots X_n$ |
|---|---|---|---|---|
| 1 | 0 | 0 | 1 | 1 |
| 2 | 1 | 1 | 0 | 0 |
| 3 | 0 | 1 | 0 | |

Ex:



$n = 3$
$T = 3$

$\text{count}(X_2 = 0 \mid X_2 = 1) = 2$
$\text{count}(X_3 = 1 \mid X_2 = 0) = 1$

* Log-likelihood:

$$\mathcal{L} = \sum_{i=1}^{n} \sum_{x} \sum_{\pi} \text{count}(X_i = x, pa_i = \pi) \log P(X_i = x_i \mid pa_i = \pi)$$

values of X $\quad$ parent configuration

* ML Estimation

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

* ML Solution (without proof):

$$P_{ML}(X_i = x \mid pa_i = \pi) = \frac{\text{count}(X_i = x, pa_i = \pi)}{\sum_{x'} \text{count}(X_i = x, pa_i = \pi)}$$

Equivalently:

$$P_{ML}(X_i = x \mid pa_i = \pi) = \frac{\text{count}(X_i = x, pa_i = \pi)}{\text{count}(pa_i = \pi)}$$

* Properties of MLE

  • Asymptotically correct: $P_{ML}(X_1, X_2, \ldots, X_n) \to P(X_1, X_2, \ldots, X_n)$ as $T \to \infty$

  • Problematic for sparse data:

  $P_{ML}(X_i = x \mid pa_i = \pi) = 0$ if $\text{count}(X_i = x \mid pa_i = \pi) = 0$

  $P_{ML}(X_i = x \mid pa_i = \pi)$ undefined if $\text{count}(pa_i = \pi) = 0$

⑤

- Other useful notation:

  Indicator function:
  $$I(x, x') = \begin{cases} 0 & \text{if } x \neq x' \\ 1 & \text{if } x = x' \end{cases}$$

  $$\text{count}(pa_i = \pi) = \sum_{t=1}^{T} I(pa_i^{(t)}, \pi)$$

  $$\text{count}(X_i = x, \; pa_i = \pi) = \sum_{t=1}^{T} I(pa_i^{(t)}, \pi) \, I(X_i^{(t)}, x)$$