Review                    7/7

* d-separation

For sets of nodes X, Y, E

when is $\begin{cases} P(Y|X,E) = P(Y|E) \\ P(X|Y,E) = P(X|E) \\ P(X,Y|E) = P(X|E)P(Y|E) \end{cases}$

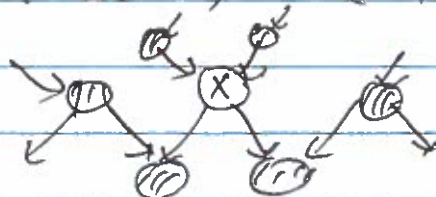* <u>True</u> if all paths from nodes in X to nodes in Y are "blocked".
A path is blocked if it has a node z that satisfies:

1)   $z \in E$    $\longrightarrow \circled{z} \longrightarrow$   intervening cause

2)   $z \in E$    $\longleftarrow \circled{z} \longrightarrow$   common cause

3)   $z \notin E$    $\longrightarrow \circled{z} \longleftarrow$   no observed common effect
     $desc(z) \notin E$

* Markov blanket $B_x$ of node X consists of parents, children and <u>spouses</u> of X.
                  other parents of
                  children

* Thm : $P(X | B_x, Y) = P(X|B_x)$   if  $Y \notin \{X, B_x\}$

\* Inference in BNs
    Query node $Q$
    Evidence nodes $E$
    How to compute $P(Q|E)$ ?

\* Polytrees
    - singly connected networks
    - polynomial time inference

\+ Loopy BNs
    - exact inference : node clustering, ....
    - approximate inference : stochastic simulation, ...

---

| 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

   (eg. learning DAGs)          (eg. learning CPTs)

- algorithms.

   non-iterative     vs      iterative

                   (loop over data many times)

  —    solution : local    vs    global   optima

\*    Maximum likelihood estimation :

- simplest form of learning in BNs
- choose ("estimate") the model ( DAG + CPTs )
    to maximize $P(\text{observed data} \mid \text{model})$

                      "likelihood"

---

| Ex: biased coin |   $P(x = \text{heads})$ |

                                         Ⓧ

    $x \in \{\text{heads, tail}\}$

    $P(x = \text{heads}) = p$

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

\* How to estimate $p$ from observed samples

   (eg $T$ coin tosses ) ?

\* IID assumption

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

$\rightarrow \{x^{(1)}, x^{(2)}, \ldots, x^{(T)}\}$   T samples

\* Probability of IID data

$$P(data) = P(X=x^{(1)}) \, P(X=x^{(2)}) \ldots P(X=x^{(T)})$$

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

\* Log - probability

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

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

log likelihood

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

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

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

$\Rightarrow N_H + N_T = T$

\* In terms of counts:

$$\ell(p) = \underbrace{N_H \log p}_{heads} + \underbrace{N_T \log (1-p)}_{tails}$$

⊁ Maximum likelihood (ML) estimation:

$$0 = \frac{dL}{dp} = \frac{N_H}{p} + \frac{N_T(-1)}{1-p} \Rightarrow N_H(1-p) + N_T p = 0$$

$$\boxed{p = \frac{N_H}{N_H + N_T} = \frac{N_H}{T}}$$  ML estimate of $p = P(heads)$
is just empirical frequency...

---

| Discrete BNs with "complete" data |

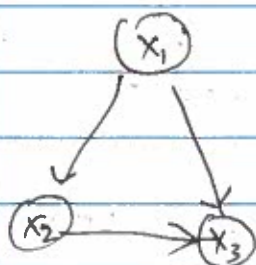\* Given : fixed DAG over discrete nodes $\{x_1, x_2 \dots x_n\}$

\* CPTs enumerate $P(X_i = x \mid \underbrace{pa(X_i)}_{parents} = \pi)$  as lookup
$\phantom{CPTs enumerate P(X_i = x)}$ of $X_i$ ↑ tables
$\phantom{CPTs enumerate P(X_i = x)}$ some configuration
$\phantom{CPTs enumerate P(X_i = x)}$ of parents.

\* Data is T complete

$\phantom{** }$ instantiations of all nodes in BN
$$\{(X_1^{(t)}, X_2^{(t)}, \dots, X_n^{(t)})\}_{t=1}^{T}$$

Ex:



n=3

| t | $X_1$ | $X_2$ | $X_3$ |
|---|---|---|---|
| 1 | 0 | 1 | 3 |
| 2 | 1 | 2 | 4 |
| 3 | 0 | 7 | 2 |
| : | 0 | : | 1 |
| : | : | : | : |
| T | 0 | 4 | 5 |

\* 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

\* IID assumption
Samples are independently, identically distributed
according to $P(X_1, X_2 \ldots X_n)$

\* Probability of (IID) data set

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

joint probability of $t^{th}$ example.

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

$$P(X_1 = x_1^{(t)}, \ldots X_n = x_n^{(t)})$$

product rule

$$= P(X_1 = x_1^{(t)}) \, P(X_2 = x_2^{(t)} | X_1 = x_1^{(t)}) \cdots P(X_n = x_n^{(t)} | X_1 = x_1^{(t)} \cdots X_{n-1} = x_{n-1}^{(t)})$$

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

$$= \prod_{i=1}^{n} P(X_i = x_i^{(t)} | pa(X_i) = pa_i^{(t)}) \quad \text{cond. ind}$$

\*    <u>Log - likelihood</u>

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

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

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

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

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

\* Let count $(X_i = x, pa_i = \Pi)$ denote # examples in table for which $X_i = x$ and $pa_i = \Pi$

Ex count $(X_2 = 2, X_1 = 0) = 1$

count $(X_2 = 1, X_1 = 0) = 2$

$\vdots$

\* log - likelihood

$$\mathcal{L} = \sum_{i=1}^{n} \sum_{x} \sum_{\Pi} \overbrace{count (X_i = x, pa_i = \Pi)}^{\text{determined by data}} \log \underbrace{P(X_i = x \mid pa_i = \Pi)}_{}$$

$\nearrow$ values that $X_i$ can assume

$\wedge$ configuration of parents of $X_i$

numbers we can choose

\* ML estimation

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

\* ML solution (without proof):

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

$\left( \text{empirical frequency of } X_i = x \text{ and } pa_i = \Pi \text{ in your data} \right)$

$$= \begin{cases} \dfrac{\text{count } (X_i = x, \ pa_i = \pi)}{\text{count } (pa_i = \pi)} & \text{when } X_i \text{ has} \\ & \qquad\qquad \text{parents} \\[2em] \dfrac{\text{count } (X_i = x)}{T} & \text{when } X_i \text{ is root} \\ & \qquad\qquad \text{node} \end{cases}$$

$T \searrow$ # examples

\*    Properties of ML estimation

− Asymptotically correct:

$$P_{ML}(X_1, X_2 \ldots X_n) \longrightarrow P(X_1 X_2 \ldots X_n) \text{ as } T \to \infty$$

− Problematic for sparse data ($T$ small)

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

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

\*    Other useful notation

Indicator function:

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

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

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

—

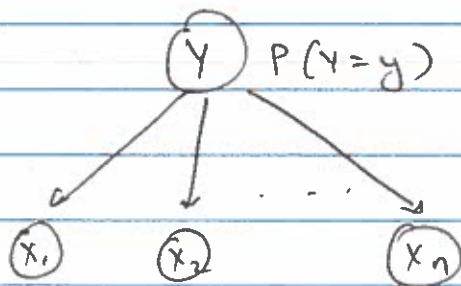Ex:  Naive Bayes model for document classification.

\* Variables

$$Y \in \{1, 2, \ldots n\}$$

eg. $1 = $ sports
$2 = $ politics

$(i = 1..n) \quad X_i \in \{0, 1\}$  does the $i^{th}$ word in the dictionary appear in document?

\*  BN = DAG + CPTs



$Y$   $P(Y=y)$

$X_1$   $X_2$   $X_n$

$P(X_i = 1 \mid Y = y)$ for $i = 1 \ldots n$

\*  How to use model for document classification.

$$P(Y = y \mid \vec{X} = \vec{x}) = \frac{P(\vec{X} = \vec{x} \mid Y = y) \, P(Y = y)}{P(\vec{X} = \vec{x})} \quad \text{Bayes rule}$$

$$= \frac{\left\{ \prod_{i=1}^{n} P(X_i = x_i \mid Y = y) \right\} P(Y = y)}{\sum_{y'} \left\{ \prod_{i=1}^{n} P(X_i = x_i \mid Y = y') \right\} P(y = y')} \quad \begin{array}{l} \text{prod rule} \\ + \text{CI} \\ \\ \text{normalization} \end{array}$$

\* How to learn a model from a large corpus of documents?

    $P_{ML}(Y=y)$    fraction of documents with topic $y$

    $P_{ML}(X_i=1 \mid Y=y)$ fraction of documents with topic $y$ that contain $i^{th}$ word in dict

\* Weaknesses of model
- "naive Bayes" assumption that words appear independently given topic

- "bag-of-words" representation ignores word ordering

<u>Ex</u>:   $\boxed{\text{Markov models of language}}$

\* Let $w_\ell$ denote the word at $\ell^{th}$ position in sentence How to model $P(w_1, w_2, \dots, w_L)$? probability of sentence with $L$ words.

\* Simplifying assumptions:

    (1) finite context/memory.    "k-gram" model

$$P(w_\ell \mid w_1, w_2, \dots w_{\ell-1}) = P(w_\ell \mid w_{\ell-1}, w_{\ell-2}, \dots w_{\ell-(k-1)})$$

                                         $k-1$ previous words

eg. $P(w_\ell \mid w_1, \dots w_{\ell-1}) = P(w_\ell \mid w_{\ell-1})$ "bi-gram" model.

(2)     position invariance (for bigram model)

$$P(w_{\ell+1} = w' \mid w_\ell = w) = P(w_{\ell+b+1} = w' \mid w_{\ell+b} = w)$$

               $b$ is any positive or negative shift.

\*    BN for bigram-model of language

$$\boxed{w_1} \longrightarrow \boxed{w_2} \longrightarrow \boxed{w_3} \rightarrow \ \cdots \ \rightarrow \boxed{w_{\ell-1}} \rightarrow \boxed{w_\ell}$$

Also : same CPT is used at all non-root
       nodes in BN.

\*    learning bigram model.

   — collect large corpus of text $\sim 10^8$ words (at least )
   — vocabulary size $V \sim 10^5$ dictionary entries

\*    Count $C_{ij} = \#$ times that word $j$ follows word $i$
    count $C_i = \#$ times that word $i$ appears in corpus

    estimate $P_{ML}(w_{\ell+1} = j \mid w_\ell = i) = \dfrac{C_{ij}}{C_i}$

\*    Note : no generalization to unseen word combinations

\* n-gram model: condition on previous n-1 words

$$P(w_\ell | w_1 \dots w_{\ell-1}) = P(w_\ell | w_{\ell-1}, \dots w_{\ell-(n-1)})$$

$$n=1 \quad \text{unigram}$$
$$n=2 \quad \text{bigram}$$
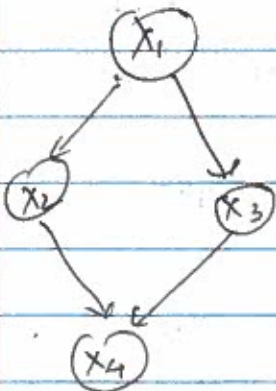$$n=3 \quad \text{trigram}$$
$$\vdots$$

n-gram counts get increasingly sparse for large n.

---

| Learning (ML estimation) from incomplete data |

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

Also: data set of T examples, but each example is a partial instantiation of nodes in BN.

Ex:



| $t$ | $X_1$ | $X_2$ | $X_3$ | $X_4$ |
|-----|-------|-------|-------|-------|
| 1 | 1 | ? | 4 | ? |
| 2 | 0 | ? | ? | 1 |
| 3 | 1 | | 5 | 3 |
| $\vdots$ | | | | |
| T-1 | 0 | | 3 | ? |
| T | 1 | ? | ? | |

\* **Goal** : estimate CPTs $P(X_i = x \mid pa_i = \pi)$ that maximize (marginal) probability of partially observed data.

\* Variables in BN

$X$ = all nodes in BN

$H$ = subset of nodes that are unobserved ("hidden")

$V$ = subset of nodes observed ("visible").

\* Log-likelihood

Assume that $T$ examples are i.i.d from joint distribution $P(X_1, X_2, \ldots, X_n)$

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

$$= \log \prod_{t=1}^{T} P(V = v^{(t)}) \quad \longleftarrow \text{visible nodes on } t^{th} \text{ example}$$

$$= \sum_{t=1}^{T} \log P(V = v^{(t)}) \quad \longleftarrow \begin{array}{c} \text{marginal probability.} \\ (\text{not joint}) \end{array}$$