## Review

* Learning in BNs
* Maximum likelihood (ML) estimation
   Estimate CPTs that maximize probability of observed data (evidence)
* Complete data (a.k.a. fully observed)
   evidence $\{x_1^{(t)}, x_2^{(t)}, \ldots, x_n^{(t)}\}_{t=1}^{T}$  T complete instantiations of nodes
$$X_1, X_2, \ldots, X_n$$

* ML estimates:
$$P_{ML}(X_i = x \mid pa_i = \pi) = \frac{Count(X_i = x, pa_i = \pi)}{\sum_{x'} Count(X_i = x', pa_i = \pi)}$$

   Equivalently:
$$P_{ML}(X_i = x \mid pa_i = \pi) = \frac{Count(X_i = x, pa_i = \pi)}{Count(pa_i = \pi)}$$

* Other notation:
   Indicator function $I(x, x') = \begin{cases} 1 & \text{if } x = x' \\ 0 & \text{otherwise} \end{cases}$
$$Count(X_i = x, pa_i = \pi) = \sum_{t=1}^{T} I(x_i^{(t)}, x)\, I(pa_i^{(t)}, \pi)$$
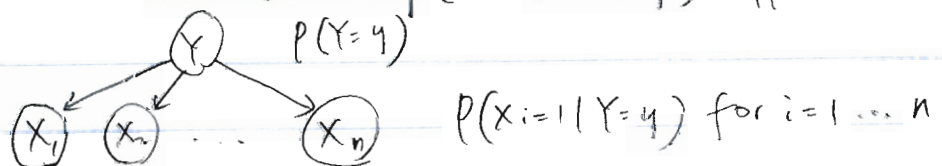
## Ex: naive Bayes model for document classification

* Variables
   $y \in \{1, 2, \ldots, m\}$ possible document topics
   $X_i \in \{0, 1\}$  does i-th word in vocabulary (dictionary) appear in documen
* BN = DAG + CPTs



$P(Y = y)$

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

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

$$P(Y=y \mid \vec{X}=\vec{x}) = \frac{\left\{ \prod_{i=1}^{n} P(X_i=x_i \mid Y=y) \right\} P(Y=y)}{\sum_{y'} \left\{ \prod_{i=1}^{n} P(X_i=x_i \mid Y=y') \right\} P(Y=y')}$$

Conditional independence

"naive Bayes" assumption

\* Strengths of model

(1) easy to estimate from a large corpus of documents

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

$P_{ML}(X_i=1 \mid Y=y)$ fraction of documents w/ topic $y$ that contain $i$-th word in vocabulary.

(2) Simplest baseline

\* Weaknesses of model

(1) naive Bayes assumption that words appear independently given topic.

(2) "bag-of-words" representation ignores word ordering

---

Ex: Markov models of language

\* Let $W_\ell$ denote word at $\ell$-th position in sentence. How to model

$P(W_1, W_2, \dots, W_{L-1}, W_L)$ probability of sentence with $L$ words $W_1, \dots, W_L$

\* Simplifying assumption

(1) finite context / memory

$$P(W_\ell \mid W_1, W_2, \dots, W_{\ell-1}) = P(W_\ell \mid W_{\ell-(k-1)}, W_{\ell-(k-2)}, \dots, W_{\ell-2}, W_{\ell-1})$$

↳ "K-gram" model

$(k-1)$ previous words

$P(W_\ell \mid W_1, W_2, \dots, W_{\ell-1}) = P(W_\ell \mid W_{\ell-1})$ "bi-gram" model

(2) position invariance

$$P(W_{\ell+1}=w' \mid W_\ell=w) = P(W_\ell=w' \mid W_{\ell-1}=w)$$

* Belief network for bigram model of language.

$$W_1 \longrightarrow W_2 \longrightarrow W_3 \longrightarrow \cdots \longrightarrow W_L$$

Same CPTs at all non-root nodes in BN.

* Learning bigram model
  * Collect large corpus of text $\sim 10^8$ words
  * 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 (followed by any word)
  estimate $P_{ML}(w_\ell = j \mid w_{\ell-1} = i) = \dfrac{c_{ij}}{c_i}$

* Note : no generalization to unseen word combinations.
* n-gram model : Condition on previous $n$ words
  $$P(w_\ell \mid w_1, \ldots, w_{\ell-1}) = P(w_\ell \mid w_{\ell-(n-1)}, \ldots, w_{\ell-1})$$
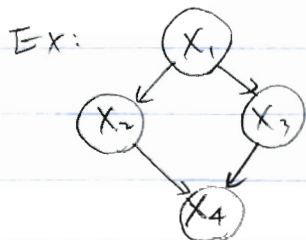
  $n=1$ unigram
  $n=2$ bigram
  $n=3$ trigram

  n-gram counts get increasingly sparse for large $n$.

[ML esitimation from incomplete data]

* Given : fixed DAG over discrete nodes $\{X_1, X_2, \ldots, X_n\}$
  Also : data set of $T$ partial instantiations of $\{X_1, X_2, \ldots, X_n\}$

EX:

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

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

\* Variables in BN

    $X =$ all nodes               $X = H \cup V$

    $H =$ hidden nodes

    $V =$ visible nodes

\* Log-likelihood

  Assume that $T$ examples are iid from joint distribution

  $P(X_1, X_2, \dots, X_n)$:

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

$$= \log \left[ \prod_{t=1}^{T} P(V = v^{(t)}) \right]$$

                   Visible nodes on $T$-th example.

$$= \sum_{t=1}^{T} \log P(V = v^{(t)})$$

$$= \sum_{t=1}^{T} \log \sum_{h} P(V = v^{(t)}, H = h) \quad \text{marginalizing over joint for } X = H \cup V$$

$$= \sum_{t=1}^{T} \log \sum_{h} \prod_{i=1}^{n} P(X_i = x \mid pa_i = \pi) \Big|_{H = h, V = v^{(t)}}$$

  \* More Complicated to optimize $\mathcal{L}$ for Incomplete data

    &minus; No "closed form" Solution.

      Alternative : iterative Solution.

\* Expectation-Maximization (EM) algorithm

  iterative procedure to maximize $\mathcal{L}(\text{data})$ for incomplete data

  in terms of CPTs.

\* Intuition — by analogy, ML estimates for complete data

$$P_{ML}(X_i = x \mid pa_i = \pi) = \frac{\text{Count}(X_i = x, pa_i = \pi)}{\text{Count}(pa_i = \pi)} = \frac{\sum_{t=1}^{T} I(x_i^{(t)}, x) \, I(pa_i^{(t)}, \pi)}{\sum_{t=1}^{T} I(pa_i^{(t)}, \pi)}$$

For incomplete data, we must "fill in" hidden values:

$$\ell_{ML}(X = x_i \mid pa_i := \pi) \longleftarrow \frac{\sum_{t=1}^{I} \ell(X_i = x, pa_i := \pi \mid V = v^{(t)})}{\sum_{t=1}^{I} \ell(pa_i := \pi \mid V^{(t)})}$$

Intuition: expected statistics ("counts") under $P(HIV)$
substitute for observed counts in complete data case.