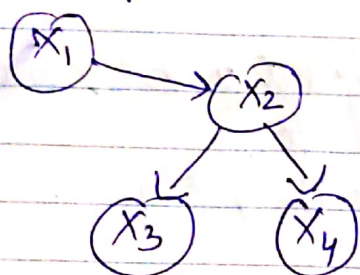→ **Learning in BNs**

Case I   Fixed DAG, complete data, lookup CPTs

Nodes $\{X_1, X_2, \ldots, X_n\}$

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

IID Data $\{(X_1^{(t)}, X_2^{(t)}, \ldots, X_n^{(t)})\}_{t=1}^{T}$

T complete instantiations of BN.



| $t$ | $X_1$ | $X_2$ | $X_3$ | $X_4$ | |
|-----|-------|-------|-------|-------|---|
| 1 | . | . | . | . | filled table |
| 2 | . | . | . | . | |
| ⋮ | . | . | . | . | |
| T | . | . | . | . | |

• **Log-likelihood of IID data**

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

$$= \sum_{i=1}^{n} \sum_{x} \sum_{\pi} \underbrace{count(X_i = x, pa_i = \pi)}_{\text{properties of data}} \cdot \underbrace{\log P(X_i = x \mid pa_i = \pi)}_{}$$

↳ Possible values of $X_i$

↳ values of parents of $X_i$

= $\pi$

↳ unknown CPTs to be estimated from data

- Write $L = \sum_{i\pi} L_{i\pi}$

  where $L_{i\pi} = \sum_x \text{count}(X_i = x, pa = \pi) \log P(X_i = x | pa_i = \pi)$

  We can independently optimize each row of each CPT in BN! [only true for complete data]

- ML Estimation

  For each node $X_i$, and for each row $\pi$ maximize $L_{i\pi}$ subject to:
  
  1) $\sum_x P(X_i = x | pa_i = \pi) = 1$

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

- Shorthand:

  Let $C_\alpha = \text{count}(X_i = \alpha, pa_i = \pi)$
  
  Let $p_\alpha = P(X_i = \alpha | pa_i = \pi)$

  How to maximize $\sum_\alpha C_\alpha \log p_\alpha$ such that

  $$p_\alpha \geq 0,$$
  $$\sum p_\alpha = 1 ?$$

  How to minimize $\sum_\alpha C_\alpha \log \frac{1}{p_\alpha}$

  equivalent to minimizing $\sum_\alpha C_\alpha \log \frac{C_\alpha}{p_\alpha}$

  $\therefore$ $C_\alpha$ are constants

Also same as minimizing $\sum_\alpha \left( \dfrac{C_\alpha}{\sum_\beta C_\beta} \right) \log \dfrac{(C_\alpha | \sum_\beta C_\beta)}{P_\alpha}$

$$\downarrow$$
KL distance

Solution: $\boxed{P_\alpha = \dfrac{C_\alpha}{\sum_\beta C_\beta}}$

## ML solution:

$$\boxed{P_{ML}(X_i = x | pa_i = \pi) = \dfrac{\text{count}(X_i = x, pa_i = \pi)}{\sum_{x'} \text{count}(X_i = x', pa_i = \pi)}}$$

- Properties

- Asymptotically correct

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

- Problematic in non-asymptotic regime (sparse data):

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

<u>Ex</u>: Markov Models of Language.

- Let $w_\ell$ denote $\ell^{th}$ word in sentence (or text)
- How to model $P(w_1, w_2, \cdots, w_L)$?
- Simplifying assumptions:

1) Finite context / history / memory:

$$P(w_\ell \mid w_1, w_2, \ldots, w_{\ell-1}) = P(w_\ell \mid w_{\ell-(n-1)}, \cdots, w_{\ell-2}, w_{\ell-1})$$

2) Position invariance:

$$P(w_\ell = w' \mid w_{\ell-(n-1)}, \cdots, w_{\ell-2}, w_{\ell-1}) =$$
$$P(w_{\ell+s} = w' \mid w_{\ell-(n-1)+s}, \cdots, w_{\ell-2+s}, w_{\ell-1+s})$$

- Markov Model

$$P(w_1, w_2, \cdots w_L) = \prod_\ell P(w_\ell \mid w_1, w_2, \cdots w_{\ell-1})$$

Product Rule

$$= \prod_\ell P(w_\ell \mid w_{\ell-(n-1)}, \cdots, w_{\ell-1}) \quad CI$$

- Models of different orders:

$n=1$ unigram $\quad (w_1) \quad (w_2) \quad (w_3) \cdots (w_n)$

$n=2$ bigram $\quad (w_1) \rightarrow (w_2) \rightarrow (w_3) - \cdots \rightarrow (w_n)$

$n=3$ trigram $\quad (w_1) \,(w_2)\,(w_3) \cdots \rightarrow (w_n)$

- Focus on bigram $(n=2)$:

  Same CPT $P(w_\ell = w' \mid w_{\ell-1} = w)$ used at each node $(\ell > 1)$

- How to learn?

  Collect large corpus of text $\sim 10^{10}$ words
  Commit to vocabulary size $\sim 10^{4-6}$

  Count $c_{ij} = \#$ times that $i^{th}$ word is followed by $j^{th}$ word in vocab.

  $c_i = \#$ times that $i^{th}$ word is followed by anything. (i.e. $c_i = \sum_j c_{ij}$)

  Estimate $P_{ML}(w_\ell = j \mid w_{\ell-1} = i) = c_{ij}/c_i$

- <u>Problems</u> with ML estimates for $n$-gram models:

  - no generalization to unseen $n$-grams.

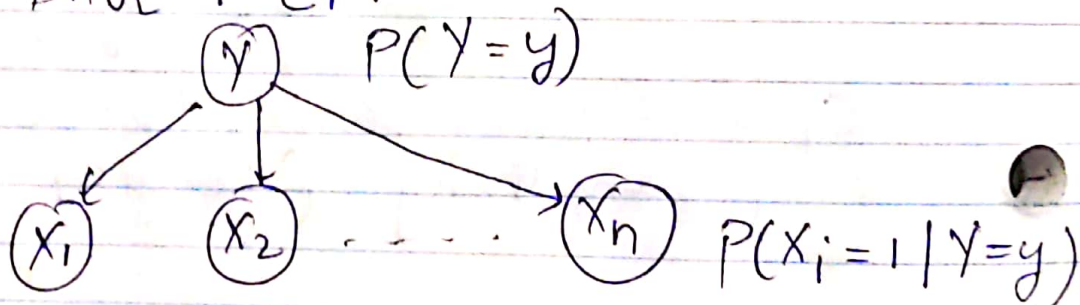  - $n$-gram counts become increasingly sparse as $n$ increases.

Ex: Naive Bayes model (for document classification)

- variables $Y \in \{1, 2, \dots m\}$ topic label
  (eg. sports, politics)

  $X_i \in \{0, 1\}$, does $i^{th}$ word appear in document?

  Used $X_i$ to represent each document as
  a fixed length $v$ vector.
  bit

- BN = DAG + CPT



  $$P(Y = y)$$

  $$P(X_i = 1 | Y = y)$$

  CPTs are unknown. How to estimate from
  data?

- How to learn?

  - Collect corpus of documents and labels
    for each document.

  - Estimate: $P_{ML}(Y = y)$ = fraction of documents
    with label $y$ in corpus.

    $P_{ML}(X_i = 1 | Y = y)$ = fraction of documents with
    label $y$ that contain
    $i^{th}$ word in dictionary

• How to classify?

$$P(Y=y \mid X_1, X_2, \ldots, X_n) = \frac{P(X_1, \ldots X_n \mid Y=y) \underset{ML}{P}(Y=y)}{P(X_1, X_2, \ldots X_n)}$$

$$= \frac{\underset{ML}{P}(Y=y) \prod_{i=1}^{n} P_{ML}(X_i \mid Y=y)}{\underset{y'}{\sum} P_{ML}(Y=y') \prod_{i=1}^{n} P_{ML}(X_i \mid Y=y')}$$
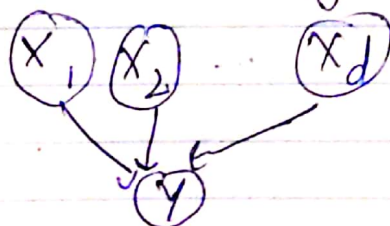
• Weaknesses:

1) "Naive" Bayes assumption that words appear independently given the topic.

2) "Bag of Words" representation ignores word ordering.

→ | Case II |  Fixed DAG, complete data,

parametrized CPT s.

(Preview) II A. linear regression



How to predict real valued $Y \in \mathbb{R}$
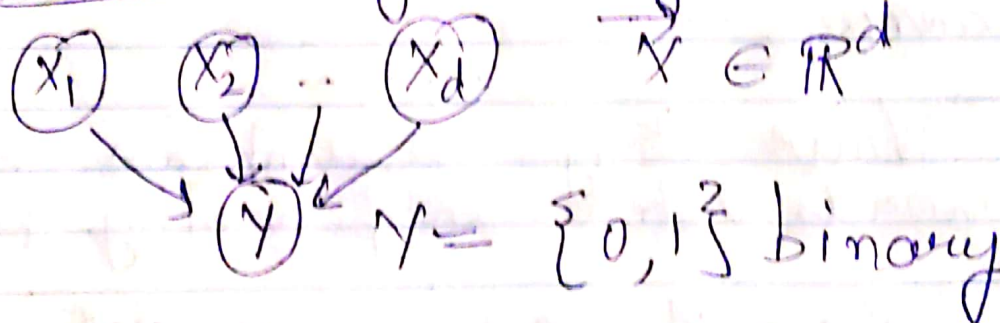from parents $\vec{X} \in \mathbb{R}^d$ ?

- Gaussian CPT

$$P(Y=y \mid \vec{X}=x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}\left(y - \sum_{i=1}^{d} w_i x_i\right)^2\right\}$$

variance

How to estimate $\vec{w}$ and $\sigma^2$ ?

$$(w_1, w_2, \ldots, w_d)$$

- Case IIB   Logistic Regression



$\vec{X} \in \mathbb{R}^d$

$Y = \{0, 1\}$ binary

How to predict binary $Y$ from a real valued $\vec{X}$

Sigmoid CPT

$$P(Y=1 \mid \vec{X}) = \sigma(\vec{w} \cdot \vec{x})$$

How to estimate $\vec{w} = (w_1, \ldots w_d)$ from data ?