

\* Learning in BNs.

\* Maximum likelihood (ML) estimation. "likelihood"

Estimate CPTs that removes probability of observed data

\* Complete data.

$x_1^{(t)}, x_2^{(t)} \dots x_n^{(t)}$  T complete instantiations of nodes  $x_1 \dots x_n$ .

\* ML Estimates.

$$P_{ML}(X_i = x | Pa_i = \pi) = \frac{\text{count}(X_i = x, Pa_i = \pi)}{\sum_{x'} \text{count}(X_i = x', Pa_i = \pi)}.$$

Equivalently.

$$P_{ML}(X_i = x | Pa_i = \pi) = \begin{cases} \frac{\text{count}(X_i = x, Pa_i = \pi)}{\text{count}(Pa_i = \pi)} & \text{for nodes w/ parents} \\ \frac{\text{count}(X_i = x)}{T} & \text{for root nodes.} \end{cases}$$

\* Other notation.

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

$$\text{Count}(X_i = x, Pa_i = \pi) = \sum_{t=1}^T I(x, x_i^{(t)}) \cdot I(\pi, Pa_i^{(t)})$$

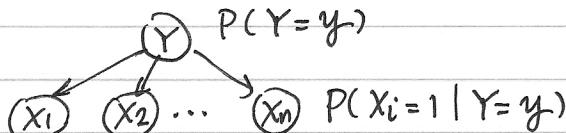
Ex: "Naive Bayes" model for document classification.

\* Variables

$Y \in \{1, 2 \dots m\}$  possible document types, topics

$X_i \in \{0, 1\}$ . = does i<sup>th</sup> word in vocabulary (dictionary) appear in document?

\* BN = DAG + CPTs



\* ML estimation of CPTs

Collect and label a large corpus of documents

$$P_{ML}(Y=y) = \text{Fraction of documents w/ topic } y$$

$P_{ML}(X_i=1 | y)$  = fraction of documents of topic  $y$   
that contains  $i^{\text{th}}$  word in vocabulary.

\* Document Classification.

$$P(Y=y | \vec{X} = \vec{x}) = \frac{P(\vec{X} = \vec{x} | Y=y) \cdot P(Y=y)}{P(\vec{X} = \vec{x})} \quad \text{Bayes rule.}$$

$$= \frac{\prod_{i=1}^n P(X_i = x_i | Y=y) \cdot P(Y=y)}{\sum_y \prod_{i=1}^n P(X_i = x_i | Y=y') \cdot P(Y=y')} \quad \begin{matrix} \text{Conditional} \\ \text{independence.} \end{matrix}$$

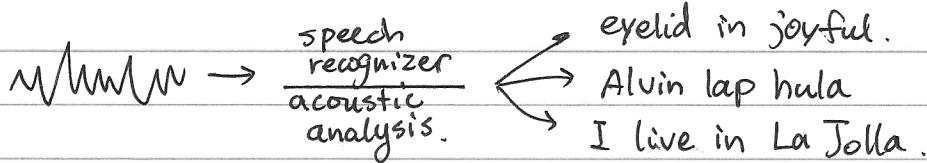
\sum \text{ sum of nominators}

\* Weakness of this "Naive Bayes" model.

- (1). Assumption that words appear independently given the topic
- (2). Documents have only one topic
- (3). "Bag of words" representation. : ignores order.

Ex: Markov model of language

why do we need language models.



\* Let  $w_l$  denote word at  $l^{\text{th}}$  position in sentence

How to model  $P(w_1, w_2 \dots w_L)$ ?

Probability of sentence with  $L$  words  $w_1, w_2 \dots w_L$ .

\* Simplifying assumptions

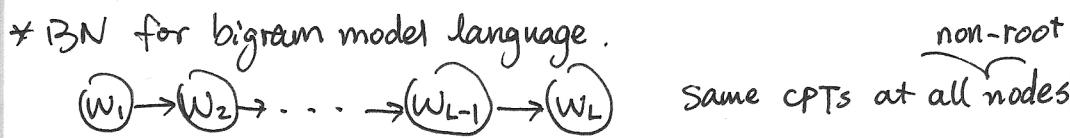
(1). finite context/memory.

$$P(w_l | \underbrace{w_1, w_2 \dots w_{l-1}}_{\text{all previous words}}) = P(w_l | \underbrace{w_{l-(k-1)}, w_{l-(k-2)}, \dots, w_{l-1}}_{k-1 \text{ previous words}}) \quad \text{"k-gram"}$$

~ special case : "bigram" model.

(2). position invariance.

$$P(w_{l+1} = w' | w_l = w) = P(w_l = w' | w_{l-1} = w)$$



\* Learning bigram model.

- collect large corpus of text  $\sim 10^8$  words
- vocabulary size  $\sqrt{\sim 10^5}$  dictionary entries
- count  $C_{ij} = \#$  times that word  $j$  follows word  $i$
- count  $C_i = \#$  times that word  $i$  appears.

$$\text{estimate } P_{ML}(W_t=j | W_{t-1}=i) = \frac{C_{ij}}{C_i}$$

\* Note: no generalization to unseen word combinations.

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

$$P(W_t | W_1, \dots, W_{t-1}) = P(W_t | W_{t-(n-1)}, \dots, W_{t-1})$$

n-gram model counts get more sparse as  $n$  increases.

### ML estimation from incomplete data

\* Given fixed graph (DAG) over discrete nodes  $\{x_1, x_2, \dots, x_n\}$ .

Also data set of  $T$  partial instantiations of  $\{x_1, x_2, \dots, x_n\}$ .

Ex:

t	$x_1$	$x_2$	$x_3$	$x_4$
1	0	?		1
2	1	2	?	1
3	0	?	?	1
4	?	?	?	0

\* Goal is to estimate CPTs  $P(x_i=x | Pa_i=\pi)$  that maximizes the.

marginal (not joint) probability of partially observed data.

(not complete)

\* Variables in BN.

$$X = \text{all nodes} \quad X = H \cup V.$$

$$H = \text{hidden nodes}$$

$$V = \text{visible nodes}$$

\* Log-likelihood. Assume  $T$  examples are i.i.d. from joint distribution

$$P(x_1, x_2, \dots, x_n)$$

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

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

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

$$\begin{aligned}
 &= \sum_{t=1}^T \log \sum_h P(V^{(t)} = v^{(t)}, H^{(t)} = h) \\
 &= \sum_{t=1}^T \log \cdot \sum_h \prod_{i=1}^n P(x_i = x_i | p_{ai} = \pi_i) \Big|_{\substack{V^{(t)} = v^{(t)} \\ H^{(t)} = h}}
 \end{aligned}$$

\* More complicated to optimize  $L$  from incomplete data.

- No "closed-form" solution.
- Iterative solution.