

Review

• Learning in BNs

• Maximum likelihood (ML) estimation

Estimate CPTs that maximize

"likelihood"

probability of observed data

• Complete data

$\left\{ (x_1^{(t)}, x_2^{(t)}, \dots, x_n^{(t)}) \right\}_{t=1}^T$ T complete instantiations of nodes X_1, \dots, X_n

• Notation

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

• ML estimates

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

$$= \begin{cases} \frac{\sum_{t=1}^T I(x_i^{(t)}, x) I(pa_i^{(t)}, \pi)}{\sum_{t=1}^T I(pa_i^{(t)}, \pi)} \\ \frac{\sum_{t=1}^T I(x_i^{(t)}, x)}{T} \end{cases}$$

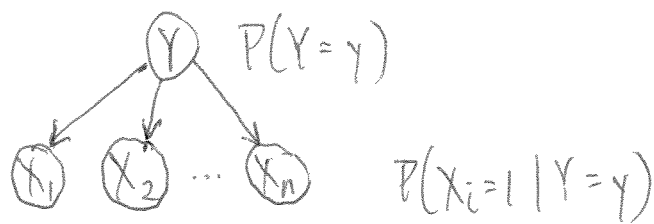
Ex: Naive Bayes model for document classification

• Variables

$Y \in \{1, 2, \dots, m\}$ possible document topics (e.g. $Y \in \{\text{spam}, \text{not spam}\}$)

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

• BN = DAG + CPTs



• ML estimation of CPTs

Collect and label a large corpus of N documents

$$P_{ML}(Y=y) \text{ ~~is the fraction of documents with topic } y~~$$

$$= \frac{\text{count}(Y=y)}{N}$$

= fraction of documents w/ topic y

$$P_{ML}(X_i=1 | Y=y) = \frac{\text{count}(X_i=1, Y=y)}{\text{count}(Y=y)}$$

= fraction of documents of topic y that contain i^{th} word in vocabulary

- Document classification

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

$$= \frac{\prod_{i=1}^n P(x_i = x_i | Y=y) P(Y=y)}{\sum_{y'} P(\vec{X}=\vec{x}, Y=y')}$$

conditional independence

$$\sum_{y'} P(\vec{X}=\vec{x}, Y=y')$$

marginalization

$$= \frac{\prod_{i=1}^n P(x_i = x_i | Y=y) P(Y=y)}{\sum_{y'} \prod_{i=1}^n P(x_i = x_i | Y=y') P(Y=y')}$$

product rule
+
cond. ind.

- Strengths of model

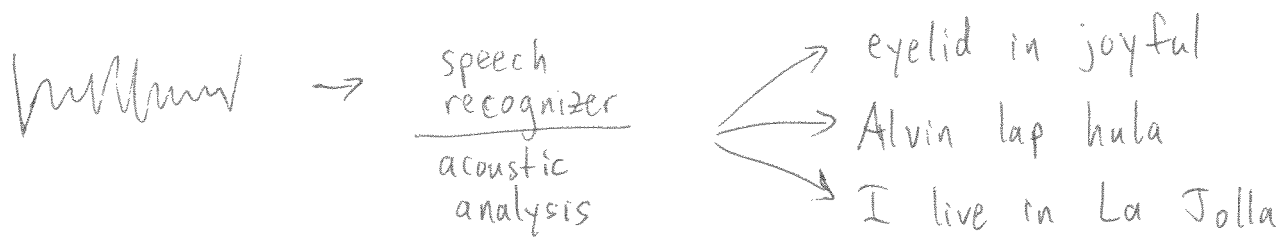
- (1) easy to estimate from a large corpus of documents
- (2) simplest baseline

- Weaknesses of model

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

Ex: Markov models of language

Why do we need language models?



- Let w_l denote word at l^{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

$$P(w_l | \underbrace{w_1, w_2, \dots, w_{l-1}}_{\text{all previous words}}) = P(w_l | \overbrace{w_{l-(n-1)}, w_{l-(n-2)}, \dots, w_{l-1}}^{n-1 \text{ previous words}})$$

"n-gram" model

special case: "bigram" model

$$P(w_l | w_1, \dots, w_{l-1}) = P(w_l | w_{l-1})$$

(2) position invariance

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

- BN for bigram model of language



- Learning bigram model

- collect large corpus of text, $\sim 10^8$ words

- vocabulary size $V \sim 10^5$ dictionary entries

- count $C_i = \#$ times word i appears

- $C_{ij} = \#$ times word j follows word i

estimate
$$P_{ML}(w_l = j | w_{l-1} = i) = \frac{C_{ij}}{C_i}$$

- Note: no generalization to unseen word combinations
(will have 0 probability)

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

- $n=1$ unigram

- $n=2$ bigram

- $n=3$ trigram

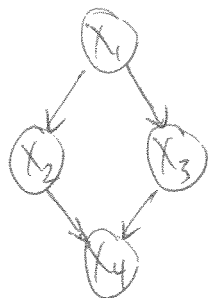
$$P(w_l | w_1, \dots, w_{l-1}) = P(w_l | w_{l-(n-1)}, \dots, w_{l-1})$$

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 | 1 |
| 2 | 1 | ? | ? | 1 |
| 3 | 0 | ? | 1 | 1 |
| 4 | ? | ? | ? | 0 |

- Goal: estimate CPTs $P(X_i = x | \text{pa}_i = \pi)$ that maximize the marginal (not joint) probability of partially observed data. (not complete)

- Variables in BN

X = all nodes

$X = H \cup V$

H = hidden nodes

V = visible nodes

- Log-likelihood: assume T examples are i.i.d. from joint distribution $P(X_1, X_2, \dots, X_n)$

$$\begin{aligned} L &= \log P(\text{data}) \\ &= \log \prod_{t=1}^T P(V^{(t)} = v^{(t)}) \\ &= \sum_{t=1}^T \log P(V^{(t)} = v^{(t)}) \end{aligned}$$

↙ visible nodes on t^{th} example

$$= \sum_{t=1}^T \log \sum_h P(V^{(t)} = v^{(t)}, H^{(t)} = h) \quad \begin{array}{l} \text{marginalizing over joint for} \\ X = HUV \\ \leftarrow \text{hidden nodes} \end{array}$$

$$= \sum_{t=1}^T \log \sum_h \prod_{i=1}^T P(X_i = x | p_{ai} = \pi) \Bigg|_{\substack{V^{(t)} = v^{(t)} \\ H^{(t)} = h}}$$

- more complicated to optimize \mathcal{L} from incomplete data
 - no "closed-form" solution
 - iterative solution

Expectation - Maximization (EM) algorithm

Iterative procedure to maximize $\mathcal{L}(\text{data})$ for incomplete data in terms of CPTs.

By analogy, ML estimates for complete data

$$P_{ML}(X_i = x | p_{ai} = \pi) = \frac{\text{count}(X_i = x, p_{ai} = \pi)}{\text{count}(p_{ai} = \pi)} = \frac{\sum_{t=1}^T I(X_i^{(t)}, x) I(p_{ai}^{(t)}, \pi)}{\sum_{t=1}^T I(p_{ai}^{(t)}, \pi)}$$

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

$$P_{ML}(X_i = x | p_{ai} = \pi) \leftarrow \frac{\sum_{t=1}^T P(X_i = x, p_{ai} = \pi | V = v^{(t)})}{\sum_{t=1}^T P(p_{ai} = \pi | V = v^{(t)})}$$

Intuition: expected statistics ("counts") under $P(H|V)$

Substitute for observed counts in complete data case