

Review Learning in BNs

1/29/15
CSE 150

* Maximum likelihood estimation (ML)

Estimate CPTs that maximize probability of observed data (evidence)

* Complete data (aka. fully observed)

Data set $\{(X_1^{(t)}, X_2^{(t)}, \dots, X_n^{(t)})\}_{t=1}^T$ is T complete instantiations of nodes X_1, X_2, \dots, X_n

t	X_1	X_2	\dots	X_n
0	0	1		0
1	1	1		0
2	0	1		1
3			...	
\vdots	\vdots	\vdots		\vdots
T	1	1		0

} data set

* ML estimates for CPTs

Nodes with parents:

$$P_{ML}(X_i = x \mid \text{pai} = \pi) = \frac{\text{count}(X_i = x, \text{pai} = \pi)}{\text{count}(\text{pai} = \pi)}$$

Root nodes:

$$P_{ML}(X_i = x) = \frac{\text{count}(X_i = x)}{T} \quad \text{where } T \text{ is \# of examples}$$

* Other notation:

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

$$\cdot \text{count}(X_i = x) = \sum_{t=1}^T I(x, x_i^{(t)}) \quad (x_i^{(t)} = i^{\text{th}} \text{ column in } t^{\text{th}} \text{ example})$$

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

* Properties of ML estimation

• Asymptotically correct

$$P_{ML}(x_1, x_2, \dots, x_n) \rightarrow P(x_1, x_2, \dots, x_n) \text{ as } T \rightarrow \infty$$

• Problematic for small #s examples (T)
("opposite limit of sparse data")

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

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

Ex: Model for document classification

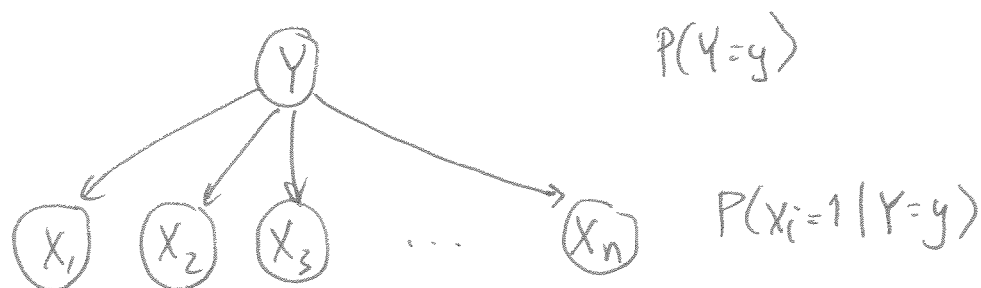
* Variables

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

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

(This converts any document to a fixed length bit vector)

* BN = DAG + CPTs



* How would you classify documents into topics with this model?

$$P(Y=y | \underbrace{\vec{X}=\vec{x}}_{\text{document's bit vector}}) = \frac{P(\vec{X}=\vec{x} | Y=y) \cdot \underbrace{P(Y=y)}_{\text{get from CPT}}}{P(\vec{X}=\vec{x})} \quad \text{Bayes rule}$$

other term in numerator:

$$P(\vec{X}=\vec{x} | Y=y) = \prod_{i=1}^n P(X_i=x_i | Y=y) \quad \text{cont. indep with d-sep case II}$$

Denominator:

$$P(\vec{X}=\vec{x}) = \sum_{y=1}^m \left\{ \prod_{i=1}^n P(X_i=x_i | Y=y) \right\} P(Y=y) \quad \text{marginalization}$$

* How to learn this model?

Estimate CPTs from a large corpus of ^Tlabeled documents

$$P_{ML}(Y=y) = \frac{\text{count}(Y=y)}{T} \quad \text{fraction of documents corresponding to topic } y \quad \text{of known topic}$$

$$P_{ML}(X_i=1 | Y=y) = \frac{\text{count}(X_i=1, Y=y)}{\text{count}(Y=y)} \quad \text{fraction of documents with topic } y \text{ that contain } i^{\text{th}} \text{ word.}$$

* Weaknesses of model

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

(ii) strong (overly strong) assumption that words are independent given topic (naive Bayes)

Ex: Markov models of language

* Let w_k denotes word at k^{th} position in sentence

How to model $P(w_1, w_2, \dots, w_{L-1}, w_L)$?

prob. of sentence with L words

* Simplifying assumptions

(i) finite context/memory

$$P(w_k | w_1, w_2, \dots, w_{k-1}) = P(w_k | \underbrace{w_{k-(k-1)}, \dots, w_{k-2}, w_{k-1}}_{\text{only condition of } (k-1) \text{ preceding words}})$$

"k-gram" model

only condition of $(k-1)$ preceding words

$$P(w_k | w_1, w_2, \dots, w_{k-1}) = P(w_k | w_{k-1}) \quad \text{is "bigram" model } (k=2)$$

(ii) position invariance

$$P(w_{k+1} = w' | w_k = w) = P(w_k = w' | w_{k-1} = w)$$

* Belief network for bigram model of language



same CPT at all non-root nodes in BN

* Learning a bigram model

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

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

* ML estimates

$$P_{ML}(w_{k+1} = j | w_k = i) = \frac{\text{count}(\text{word } i \text{ is followed by word } j)}{\text{count}(\text{word } i \text{ occurs})} = \frac{C_{ij}}{C_i}$$

parent i co-occurs with word j

* Note: no generalization to unseen word combinations

* k -gram model : conditioning on $(k-1)$ previous words

$k=1$ unigram model

$k=2$ bigram "

$k=3$ trigram "

version: more powerful model
but also more prob. where count = 0.
as k increases

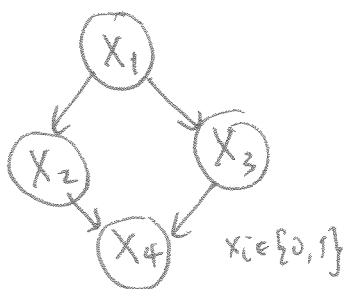
(ngrams.google.com)

Material Up To Quiz 1

Learning from incomplete data

* Given some fixed DAG over discrete nodes $\{X_1, X_2, \dots, X_n\}$
Also given 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	1	0	?	1
2	0	?	?	1
3	1	0	?	1
\vdots	\vdots	\vdots	\vdots	\vdots
T	0	1	?	1

$X_i \in \{0, 1\}$

* Goal: estimate CPTs $P(X_i = x | \text{pa}_i = \pi)$ that maximize
marginal prob. of partially observed data

(vs. before when we maximize joint prob of complete 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 $P(X_1, X_2, \dots, X_n)$ of BN

$$\begin{aligned}\mathcal{L} &= \log P(\text{DATA}) \\ &= \log \prod_{t=1}^T P(V = v^{(t)}) \quad \text{marginal prob. of visible nodes in } t^{\text{th}} \text{ example} \\ &= \sum_{t=1}^T \log P(V = v^{(t)}) \\ &= \sum_{t=1}^T \log \left[\underbrace{\sum_h P(V = v^{(t)}, H = h)}_{\text{joint}} \right] \quad \text{marginalization}\end{aligned}$$

... much more complicated to optimize b/c sum inside log.