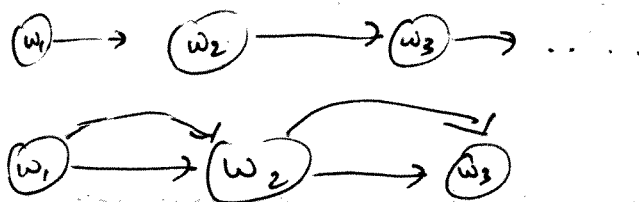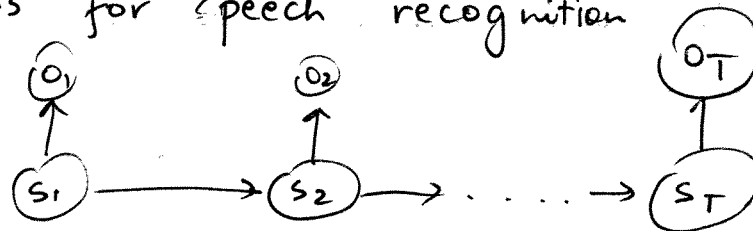Compact expressions of complex worlds

1) noisy - OR

2) naive Bayes

3) Markov models of language

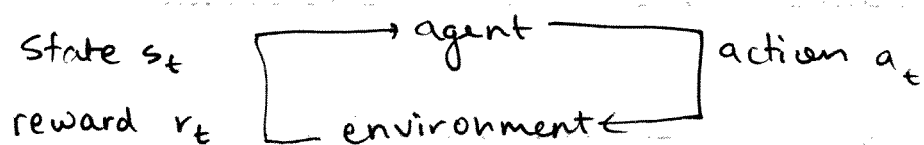$$W_\ell = \ell^{th} \text{ word in sentence}$$

$$w_1 \longrightarrow w_2 \longrightarrow w_3 \longrightarrow \dots$$

$$w_1 \longrightarrow w_2 \longrightarrow w_3$$

4) HMMs for speech recognition

$$S_1 \longrightarrow S_2 \longrightarrow \dots \longrightarrow S_T$$
with observations $O_1, O_2, O_T$

5) MDPs for planning

State $s_t$       agent       action $a_t$
reward $r_t$       environment

Efficient algorithms

1) Conditional independence tests via d-separation

(i)   $\longrightarrow \boxed{\hspace{-2pt}/\!/\!/} \longrightarrow$

(ii)  $\longleftarrow \boxed{\hspace{-2pt}/\!/\!/} \longrightarrow$

(iii) $\longrightarrow \bigcirc \longleftarrow$

2) polytree algorithm for inference



3) EM algorithm for ML estimation

update $P(X_i = x \mid pa_i = \pi) \longleftarrow$

$$\frac{\sum_t P(X_i = x, \, pa_i = \pi \mid V^{(t)})}{\sum_t P(pa_i = \pi \mid V^{(t)})}$$

guarantee = monotonic convergence

$$\mathcal{L} = \sum_t \log P(V^{(t)})$$

never decreases,

generally increases at each iteration.

4) HMM'S

- computing likelihood $P(O_1, O_2 \cdots O_T)$
- decode argmax $P(S_1, S_2 \cdots S_T \mid O_1, O_2 \cdots O_T)$
- belief updating    — learning

# 5)    Algorithms in MDPs

- Policy iteration

$$\pi_0 \xrightarrow{\text{evaluate}} \begin{array}{c} V^{\pi_0}(s) \\ Q^{\pi_0}(s,a) \end{array} \xrightarrow[\text{greedy}]{\text{improve}} \pi_1 \rightarrow \cdots$$
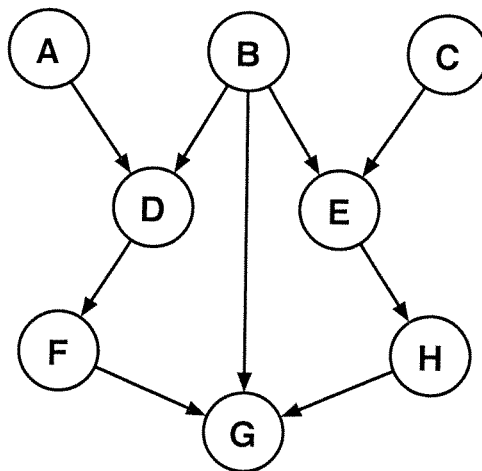
- Value iteration

$$V_{K+1}(s) \leftarrow R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) V_K(s')$$

$Q^\pi(s,a) =$ expected return if take action $a$ at $t=0$, in state $s$, then follow policy $\pi$

$$= E^\pi \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) \mid s_0 = s, a_0 = a \right]$$

## 2. Conditional independence (10 pts)

For the belief network shown below, indicate whether the following statements of conditional independence are **true (T)** or **false (F)**.



| | | | |
|---|---|---|---|
| *true* | $P(A)$ | $=$ | $P(A\|E)$ |
| *false* | $P(A\|C,G)$ | $=$ | $P(A\|G)$ |
| _____ | $P(A\|D,F,G)$ | $=$ | $P(A\|D,G)$ |
| *true* (*) | $P(A,B,C\|D,E)$ | $=$ | $P(A\|B,D)\,P(B\|C,D,E)\,P(C\|D,E)$ |
| _____ | $P(D,G\|F)$ | $=$ | $P(D\|F)\,P(G\|F)$ |
| _____ | $P(D,E\|B)$ | $=$ | $P(D\|B)\,P(E\|B)$ |
| _____ | $P(F\|C)$ | $=$ | $P(F)$ |
| _____ | $P(F\|C,E)$ | $=$ | $P(F\|E)$ |
| _____ | $P(F,H\|B)$ | $=$ | $P(F\|B)\,P(H\|B)$ |
| _____ | $P(F,H\|B,G)$ | $=$ | $P(F\|B,G)\,P(H\|B,G)$ |

(*) $P(A,B,C|D,E) = P(C|D,E)\,P(B|C,D,E)\,\underbrace{P(A|B,C,D,E)}$

$$\| \ ?$$

$$P(A|B,D)$$

3

# 3. Polytree inference



Inference in a polytree scales linearly in the number of nodes and the sizes of their conditional probability tables (CPTs). For the belief network shown above, consider how to *efficiently* compute the posterior probability $P(E|C, F)$. This can be done in four consecutive steps in which the later steps rely on the results from earlier ones.

Complete the procedure below for this inference by showing how to compute the necessary result at each step. For full credit, make each step as efficient as possible. Your answers should be expressed in terms of the CPTs of the belief network as well as the results of previous steps. **Hint:** *at each step, you'll want to exploit what you've just computed in the previous one.*
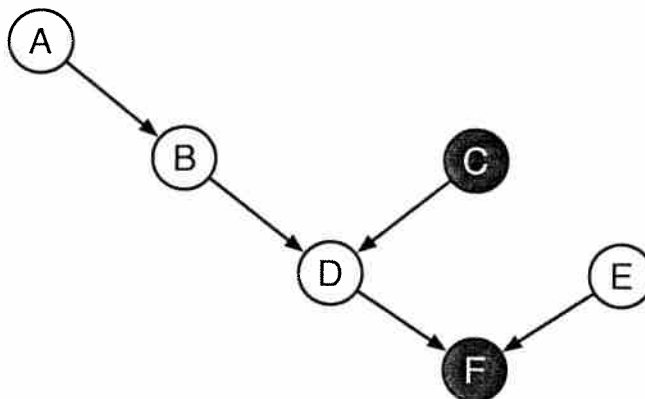
(a) **Marginal** (2 pts)

$$P(B) \;=\; \sum_a P(A=a, B) \qquad \text{Marginalization (M)}$$

$$\;=\; \sum_a P(A=a)\, P(B|A=a) \qquad \text{product rule (PR)}$$

(b) **Conditional** (3 pts)

$$P(D|C) \;=\; \sum_b P(B=b, D\,|\,C) \qquad (M)$$

$$\;=\; \sum_b P(B=b\,|\,\emptyset)\, P(D|B=b, C) \qquad (PR)\,(CI)$$

$$\;=\; \sum_b P(B=b)\, P(D|B=b, C) \qquad /\!/$$

8

## 3. Polytree inference (con't)



(c) **Conditional** (3 pts)

$P(F|C, E)$

Marginalization over D
prod rule
CI

$$= \sum_d \underbrace{P(D=d \mid C)}_{(b)} \underbrace{P(F \mid D=d, E)}_{CPT}$$

(d) **Posterior** (5 pts)

$P(E|C, F)$

$$= \frac{\overbrace{P(F|E, C)}^{(c)} P(E|C)}{P(F|C)} \qquad \text{Bayes rule}$$

$$= \frac{P(F|E,C) \; P(E)}{P(F|C)} \qquad CI$$

$$= \cancel{P(E)} \; \frac{P(F|E,C) \, P(E)}{\sum_e P(F|E=e, C) \, P(E=e)} \qquad \text{normalization}$$
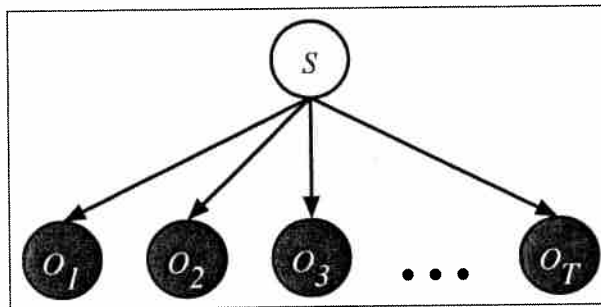
9

# 3. Belief updating

## (a) Naive Bayes model (3 pts)

Consider the belief network shown below for the discrete random variables $S \in \{1, 2, \ldots, n\}$ and $O_t \in \{1, 2, \ldots, m\}$. Also, let the CPTs of the network be parameterized by:

$$\pi_i = P(S=i),$$
$$b_{ik} = P(O_t=k|S=i),$$

where the same CPT is used for all the child nodes. In terms of these parameters, show how to compute the posterior probability $P(S=i|o_1, o_2, \ldots, o_T)$. You may use the notation $b_i(o_t)$ as shorthand for $P(o_t|S=i)$.
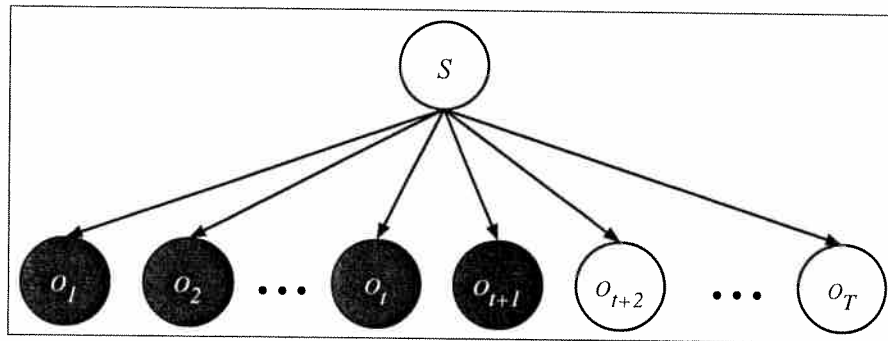


$P(s=i \mid o_1, o_2 \cdots o_T)$

$= \dfrac{P(o_1, o_2 \cdots o_T \mid s=i) \, P(s=i)}{P(o_1 \cdots o_T)}$  Bayes

$= \dfrac{P(s=i) \times \left( P(o_1 \mid s=i) \, P(o_2 \mid s=i) \cdots P(o_T \mid s=i) \right)}{P(o_1 \cdots o_T)}$  Product rule, CI

$= \dfrac{\pi_i \displaystyle\prod_{t=1}^{T} P(o_t \mid s=i)}{\displaystyle\sum_{j=1}^{n} \pi_j \prod_{t=1}^{T} P(o_t \mid s=j)}$  normalization

14

## (b) Belief updating (5 pts)

Consider the same belief network as before. In this problem, you will derive an incremental update rule for the scenario where the observations arrive one at a time. Let

$$q_{it} = P(S{=}i|o_1, o_2, \ldots, o_t)$$

denote the posterior probability over the hidden variable $S$ based on the first $t$ observations. Derive a recursion for these probabilities at time $t{+}1$ in terms of those at time $t$. (Your answer should also involve the parameters of the belief network and the observation $o_{t+1}$ at time $t{+}1$.)



**Note**: the incremental update should require only $O(n)$ operations, **not** $O(nt)$ or $O(nT)$ operations as presumably did your answer to part (a).
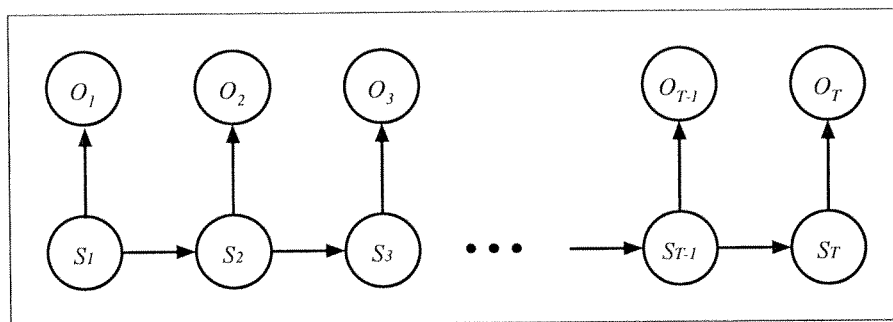
$$q_{i,t} = P\left(s=i \mid o_1, \ldots o_t\right)$$

$$q_{i,t+1} = P\left(s=i \mid o_1, o_2, \ldots, o_{t+1}\right)$$

$$= \frac{P\left(o_{t+1} \mid s=i, o_1, \ldots o_t\right) P\left(s=i \mid o_1, o_2 \cdots o_t\right)}{P\left(o_{t+1} \mid o_1 o_2 \cdots o_t\right)} \quad \text{Bayes rule}$$

$$= \frac{P\left(o_{t+1} \mid s=i\right) P\left(s=i \mid o_1, \cdots o_t\right)}{P\left(o_{t+1} \mid o_1, \ldots o_t\right)} \quad \text{CI}$$

$$= \frac{b_i\left(o_{t+1}\right) q_{it}}{\sum_j b_j\left(o_{t+1}\right) q_{jt}} \qquad \begin{array}{l} \text{substitution} \\[4pt] \text{normalization} \end{array}$$

15

### (c) Comparison to HMM (2 pts)

Now consider a discrete HMM where as usual $S_t \in \{1, 2, ..., n\}$ denotes the hidden state at time $t$, $O_t \in \{1, 2, ..., m\}$ denotes the observation at time $t$, and the parameters

$$
\begin{aligned}
a_{ij} &= P(S_{t+1}=j|S_t=i), \\
b_{ik} &= P(O_t=k|S_t=i),
\end{aligned}
$$

denote the transition and emission matrices. Also let $q_{it} = P(S_t = i|o_1, o_2, \ldots, o_t)$ denote the posterior over the hidden state $S_t$ at time $t$ after seeing $t$ observations. In an HMM, the equation for belief updating is given by:

$$
q_{j,t+1} = \frac{\sum_i q_{it} a_{ij} b_j(o_{t+1})}{\sum_{i'j'} q_{i't} a_{i'j'} b_{j'}(o_{t+1})}.
$$

Suppose we consider the trivial HMM in which the hidden state *never* changes. This can be enforced by setting the transition matrix $a_{ij}$ equal to the identity matrix: that is,
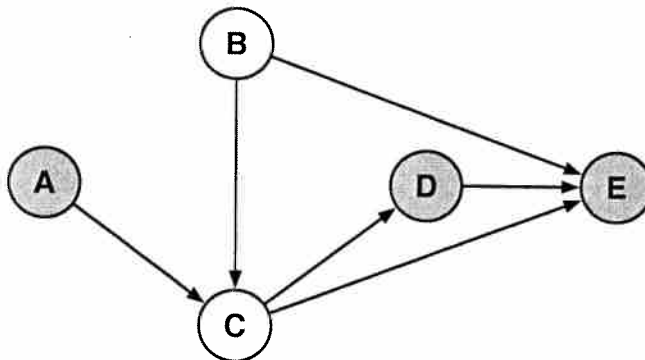
$$
a_{ij} = I(i, j),
$$

where $I(i, j)$ is the usual indicator function. For this special case, does the above formula for belief updating reduce to your answer in part (b)? If yes, show that it does; if not, explain why it doesn't.

$$
q_{j,t+1} = \frac{\sum_i q_{it} a_{ij} b_j(o_{t+1})}{\sum_{i',j'} q_{i't} a_{i'j} b_{j'}(o_{t+1})}
$$

$$
= \frac{\sum_{i=1}^{n} q_{it} I(i,j) b_j(o_{t+1})}{\lfloor \qquad \qquad \rfloor}
$$

16

$$
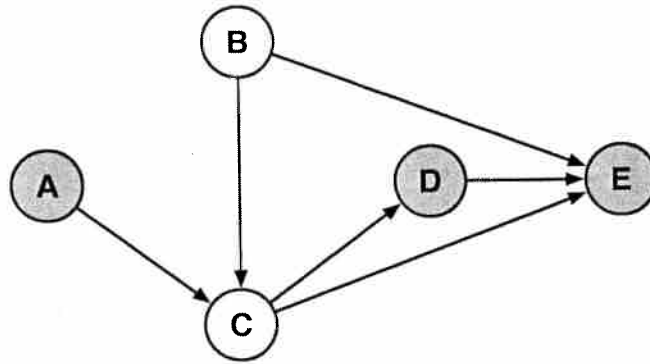= \frac{q_{jt} b_j(o_{t+1})}{\sum_{j'} q_{j't} b_{j'}(o_{t+1})}
$$

# 1. EM algorithm

Consider the belief network, shown below, with conditional probability tables $P(A)$, $P(B)$, $P(C|A, B)$, $P(D|C)$, and $P(E|B, C, D)$. Suppose that values of the nodes $B$ and $C$ are unobserved (hidden), while the values of the shaded nodes $A$, $D$, and $E$ are observed.



(a) **Posterior probability** (4 pts)

Show how to compute the posterior probability $P(B = b, C = c|A = a, D = d, E = e)$ in terms of the belief network's conditional probability tables. You may use shorthand notation to simplify your answer: i.e., $P(b, c|a, d, e) = P(B = b, C = c|A = a, D = d, E = e)$.

**(b) Posterior probability** (2 pts)

Compute the posterior probabilities $P(b|a, d, e)$ and $P(c|a, d, e)$ in terms of your answer from part (a); in other words, for this problem you may assume that the posterior probability $P(b, c|a, d, e)$ is given.
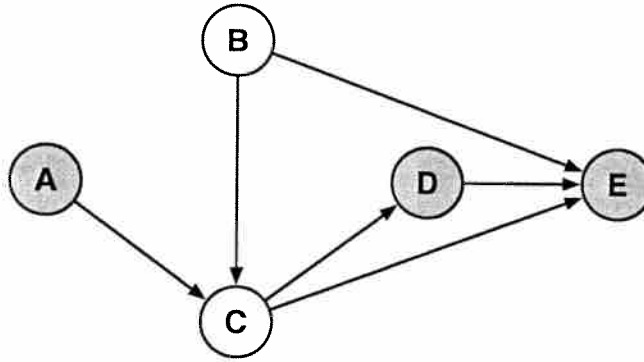
**(c) Log-likelihood** (3 pts)

Consider a data set of $T$ partially labeled examples $\{a_t, d_t, e_t\}_{t=1}^T$ over the nodes $A$, $D$, and $E$ of the belief network. The log-likelihood of the data is given by:

$$\mathcal{L} = \sum_{t=1}^{T} \log P(A=a_t, D=d_t, E=e_t)$$

Compute this log-likelihood in terms of the CPTs of the belief network.

Marginalization

$$= \sum_t \log \sum_{b,c} P\left(a_t, B=b, C=c, d_t, e_t\right)$$

$$= \sum_t \log \sum_{b,c} P(a_t) P(b) P(c|a_t, b) P(d_t|c) P(e_t|b, d_t)\Big]$$

PR & CI

3

(d) **EM algorithm** (8 pts)

Consider the EM algorithm that updates the CPTs to maximize the log-likelihood of the data set in part (c). Complete the numerator and denominator in the expressions for the EM updates shown below and on the following pages.
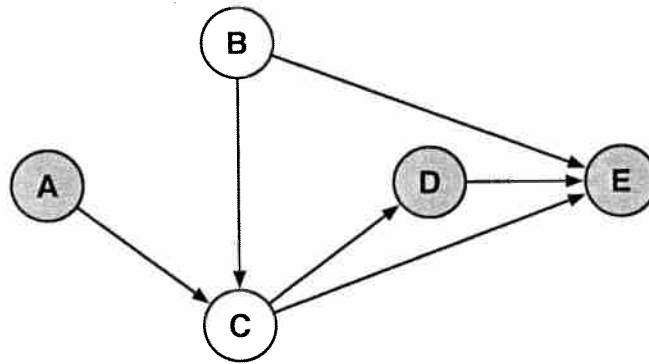
*Suggested notation.* Use shorthand such as $P(b, c|a_t, d_t, e_t)$ for the posterior probabilities computed in parts (a) and (b). Also, use indicator functions such as $I(a, a_t)$, where:

$$I(a, a_t) = \begin{cases} 1 \text{ if } a = a_t, \\ 0 \text{ if } a \neq a_t. \end{cases}$$

Simplify your answers as much as possible, and put your *final* answers in the spaces provided below.

$$P(B=b) \leftarrow \frac{\sum_{t=1}^{T} P(b|a_t, d_t, e_t)}{T}$$

$$P(C=c|A=a, B=b) \leftarrow \frac{\sum_{t=1}^{T} I(a, a_t) P(b, c|a_t, d_t, e_t)}{\sum_{t=1}^{T} I(a, a_t) P(b|a_t, d_t, e_t)}$$

4

(d) **EM algorithm** (continued)

$$P(D{=}d|C{=}c) \quad \leftarrow \quad \underline{\hspace{8cm}}$$

$$P(E{=}e|B{=}b, C{=}c, D{=}d) \quad \leftarrow \quad \frac{\sum_t I(d, d_t)\, I(e, e_t)\, P(b, c \,|\, a_t, d_t, e_t)}{\sum_t I(d, d_t)\, P(c, b \,|\, a_t, d_t, e_t)}$$

(e) **Multiple choice** (3 pts)

Consider the EM algorithm for belief networks of discrete nodes with arbitrarily specified CPTs. Indicate the *best* answer to each statement in the space provided.

_____ **In practice, the EM algorithm for these belief networks will generally:**

(a) converge to a *global maximum* of the log-likelihood.

(b) converge to a *local maximum* of the log-likelihood.

(c) converge to a *global minimum* of the log-likelihood.

(d) converge to a *local minimum* of the log-likelihood.

(e) not converge.

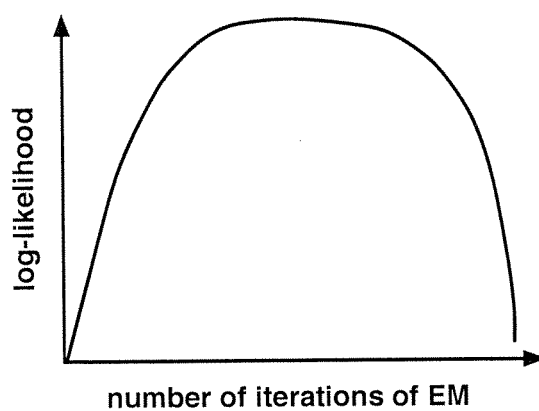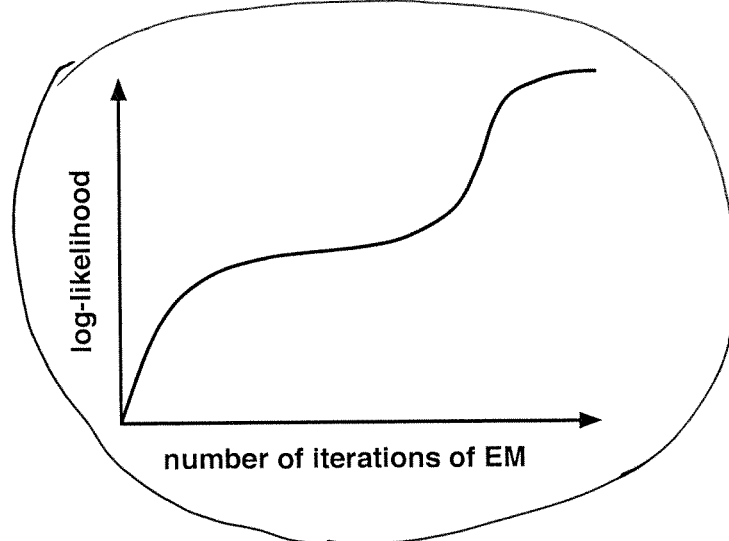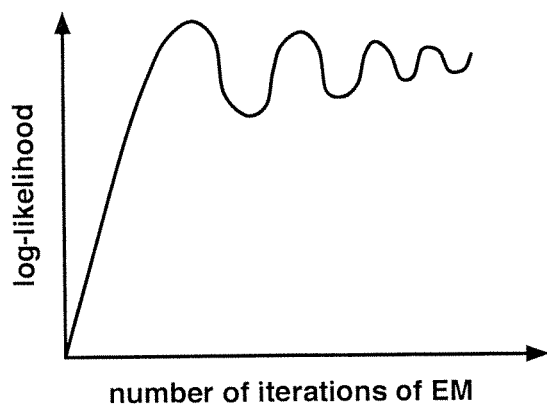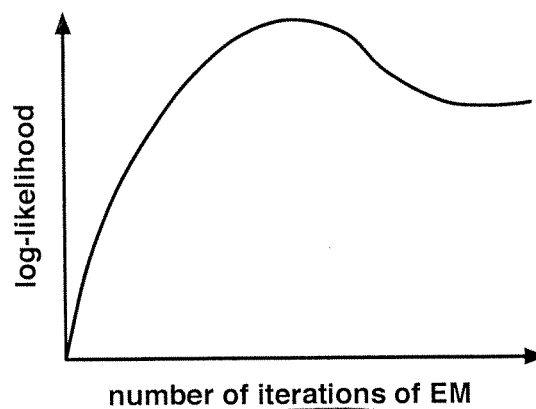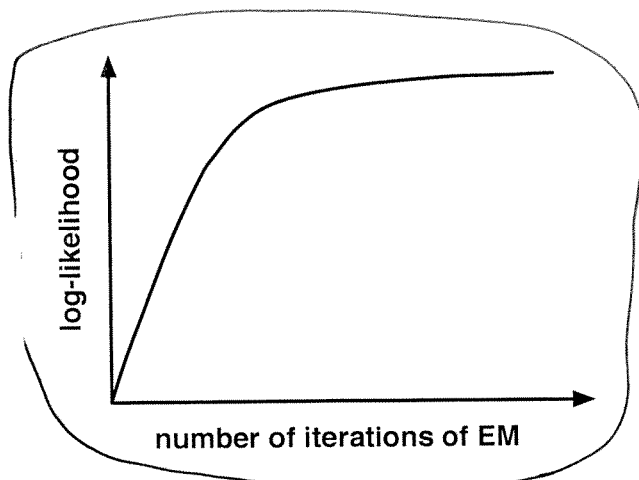_____ **The EM algorithm for these belief networks is guaranteed to converge:**

(a) within $T$ iterations, where $T$ is the number of examples.

(b) only when the parameters are initialized by a domain expert.

(c) as long as there are no duplicate examples in the data set.

(d) monotonically, never decreasing the log-likelihood at each iteration.

(e) none of the above (unless the belief network happens to be a polytree).

_____ **In general each iteration of the EM algorithm scales:**

(a) linearly in the number of examples, $T$.

(b) linearly in the number of edges of the belief network.

(c) inversely in the number of observed nodes.

(d) exponentially in the number of hidden nodes.

(e) all of the above.

(f) none of the above.
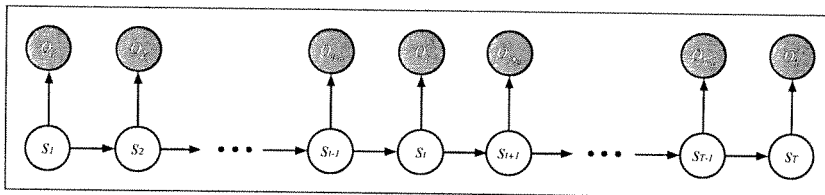
6

## (e) Convergence (5 pts)

The EM algorithm generally converges to a local maximum of the log-likelihood. Circle the plots below that might be obtained from a correct implementation of the EM algorithm.

# 2. Hidden Markov modeling

Consider a discrete hidden Markov model (HMM) with the belief network shown below. Let $S_t \in \{1, 2, ..., n\}$ denote the hidden state of the HMM at time $t$, and let $O_t \in \{1, 2, ..., m\}$ denote the observation at time $t$. As in class, let:

$$\begin{aligned}
\pi_i &= P(S_1 = i), \\
a_{ij} &= P(S_{t+1} = j | S_t = i), \\
b_{ik} &= P(O_t = k | S_t = i),
\end{aligned}$$



denote the parameters of the HMM. You may also use $b_i(k)$ to denote the matrix element $b_{ik}$.

(a) **Backward algorithm** (4 pts)

Consider the probabilities $\beta_{it} = P(o_{t+1}, o_{t+2}, \ldots, o_T | S_t = i)$. Derive an efficient recursion to compute these probabilities at time $t < T$ from those at time $t+1$ and the transition and emission matrices of the HMM. Justify briefly each step in your derivation.

$$\beta_{it} = P(o_{t+1}, o_{t+2}, \ldots, o_T | S_t = i)$$

$$= \sum_j P(s_{t+1} = j, o_{t+1}, \ldots o_T | S_t = i) \quad \text{marginalization}$$

$$= \sum_j P(s_{t+1} = j | S_t = i) \, P(o_{t+1} | \cancel{S_t = i}, S_{t+1} = j) \, P(o_{t+2} \ldots o_T | \cancel{S_t = i}, S_{t+1} = j, \cancel{o_{t+1}}) \quad \overset{PR \ \& \ CI}{}$$

$$= \sum_j P(s_{t+1} = j | S_t = i) \, P(o_{t+1} | s_{t+1} = j) \, P(o_{t+2} \ldots o_T | s_{t+1} = j)$$

$$= \sum_j a_{ij} \, b_j(o_{t+1}) \, \beta_{j, t+1}$$

6

(d) **Viterbi algorithm** (6 pts)

For a particular sequence of observations $(o_1, o_2, \ldots, o_T)$, suppose that recursive algorithms are implemented to compute the $n \times T$ matrices (from lecture) with elements:

$$\ell_{it}^* = \max_{s_1,\ldots,s_{t-1}} \left[ \log P(s_1, \ldots, s_{t-1}, s_t = i, o_1, \ldots, o_t) \right]$$

$$\Phi_{j,t+1} = \operatorname*{argmax}_i \left[ \ell_{it}^* + \log a_{ij} \right]$$

(i) From the matrix elements shown above, which you may assume to be given, show how to compute the most likely hidden state sequence

$$\{s_1^*, s_2^*, \ldots, s_T^*\} = \operatorname*{argmax}_{s_1,\ldots,s_T} P(s_1, s_2, \ldots, s_T | o_1, o_2, \ldots, o_T).$$

In particular, do this by completing the following pseudocode for the Viterbi algorithm (making use of the matrices defined above):

$$s_T^* = \operatorname*{argmax}_i \left[ \ell_{iT}^* \right]$$
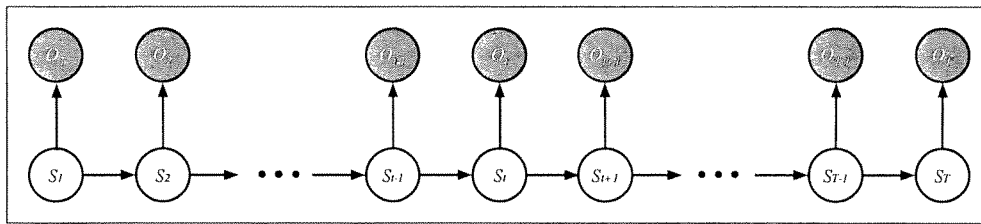
for $t = T-1$ to 1

$$s_t^* = \Phi\left( s_{t+1}^*, t+1 \right)$$

(ii) Shown below are the matrix elements $\Phi_{it}$ (for $i = 1$ to $n$ and $t = 1$ to T) for a particular observation sequence $\{o_1, \ldots, o_T\}$ of an $n$-state HMM with $n = 5$ and $T = 6$. Given that $s_T = 4$, use the matrix elements to deduce the missing values of the most likely state sequence:

| 4 | 4 | 3 | 1 | 3 | 4 |
|---|---|---|---|---|---|
| $s_1$ | $s_2$ | $s_3$ | $s_4$ | $s_5$ | $s_6$ |

$$\Phi = \begin{pmatrix} * & 1 & 2 & \boxed{3} & 4 & 5 \\ * & 2 & 3 & 4 & 5 & 1 \\ * & 3 & \boxed{4} & 5 & \boxed{1} & 2 \\ * & \boxed{4} & 5 & 1 & 2 & \boxed{3} \\ * & 5 & 1 & 2 & 3 & 4 \end{pmatrix}$$

**Note:** the asterisk symbols ($*$) indicate undefined elements at time $t = 1$.

(d*) **Most likely hidden state** (2 pts)

As before, suppose that for a particular sequence of observations, the forward-backward algorithm in HMMs is used to compute the probabilities:

$$\alpha_{it} = P(o_1, o_2, \ldots, o_t, S_t = i),$$
$$\beta_{it} = P(o_{t+1}, o_{t+2}, \ldots, o_T | S_t = i).$$

In terms of these probabilities, which you may assume to be given, show how to compute the most likely value of the (single) hidden state at time $t$:

$$\tau_t^* = \operatorname*{argmax}_i \left[ P(S_t = i | o_1, o_2, \ldots, o_T) \right].$$

Show your work for full credit, briefly justifying each step in your derivation.

$$= \operatorname*{argmax}_i \left[ \frac{P(S_t = i, o_1, o_2, \ldots o_T)}{P(o_1 \ldots o_T)} \right] \quad \text{prod rule}$$

$$= \operatorname*{argmax}_i \left[ P(S_t = i, o_1, \ldots o_T) \right] \quad \text{den. ind of } i$$

$$= \operatorname*{argmax}_i \left[ P(S_t = i, o_1 \ldots o_t) P(o_{t+1} \ldots o_T | S_t = i, o_T \ldots o_t) \right]$$
$$\text{prod rule, CI}$$

$$= \operatorname*{argmax}_i \left[ \alpha_{it} \beta_{it} \right]$$

9

(e*) **Compare and contrast** *(2 pts)*

The Viterbi algorithm computes the most likely *sequence* of hidden states for a particular sequence of observations:

$$\{s_1^*, s_2^*, \ldots, s_T^*\} = \underset{\{s_1, s_2, \ldots, s_T\}}{\operatorname{argmax}} \left[ P(s_1, s_2, \ldots, s_T | o_1, o_2, \ldots, o_T) \right]$$

Consider how these collectively optimal states $s_t^*$ differ (if at all) from the individually optimal states $\tau_t^*$ defined in part (d) of this problem:

$$\tau_t^* = \underset{i}{\operatorname{argmax}} \left[ P(S_t = i | o_1, o_2, \ldots, o_T) \right].$$

Answer the following [yes/no] questions:

**No**      Is it always true that $\tau_t^* = s_t^*$ for all $t$?

**no**      Is it always true that $P(\tau_1^*, \tau_2^*, \ldots, \tau_T^*) > 0$?

**Yes**      Is it possible that $\tau_t^* = s_t^*$ for all $t$?

**never !!**      Is it possible that $P(\tau_1^*, \ldots, \tau_T^* | o_1, \ldots, o_T) > \underbrace{P(s_1^*, \ldots, s_T^* | o_1, \ldots, o_T)}$?

$S_t^*$ defined to maximize this

10

1. Suppose $X, Y$, and $Z$ are random variables. Which of the following expressions can be used to compute $P(X|Y, Z)$?

   (a) $P(X, Y, Z)/P(Y, Z)$

   (b) $P(Y|X, Z)P(X|Z)/P(Y|Z)$

   (c) $P(Y, Z|X)P(X)/P(Y, Z)$

   (d) some of the above

   (e) all of the above

2. In a polytree belief network,

   (a) each node has no more than one parent.

   (b) there is at most one undirected path between any two nodes.

   (c) the graph is a disjoint union of trees.

   (d) each loop in the graph contains exactly one root node.

3. A naive Bayes model for classification of multivariate data assumes that

   (a) the features of the data are marginally independent.

   (b) the features of the data are conditionally independent given the label.

   (c) the features of the data are uniformly distributed within the unit hypercube.

   (d) there are fewer class labels than features.

4. Consider a discrete hidden Markov model with $n$ hidden states and $m$ observations. Given an observation sequence of length $T$, the Viterbi algorithm computes the most likely hidden state sequence in time

   (a) $O(n^2 m^2 T)$

   (b) $O(m^2 T)$

   (c) $O(n^2 T)$

   (d) $O(n^T)$

5. Let $W_k$ represent the $k$th word in a sentence. Which of the following is true in a trigram model of natural language?

   (a) $P(W_k) = P(W_k|W_{k-1}, W_{k-2})$

   (b) $P(W_{k-1}, W_k, W_{k+1}) = P(W_{k-1})P(W_k)P(W_{k+1})$

   (c) $P(W_k|W_{k-1}, W_{k-2}) = P(W_k|W_1, W_2, \ldots, W_{k-1})$

   (d) $P(W_k|W_{k-1}, W_{k-2}) = P(W_k|W_1, W_2, \ldots, W_{k-1}, W_{k+1}, W_{k+2}, \ldots, W_{\text{final}})$

1