(1)

\*    Belief network = DAG + CPTs

$$P(X_1, X_2, \ldots X_n) = \prod_{i=1}^{n} P(X_i | X_1, X_2 \ldots, X_{i-1})$$

$$= \prod_{i=1}^{n} \underbrace{P(X_i | pa_i)}_{CPTs}$$

+    Learning from complete data

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

$$P_{ML}(X_i = x | pa_i = \pi) = \frac{count(X_i = x, \; pa_i = \pi)}{count(pa_i = \pi)}$$

$$= \frac{\sum_t I(X_i^{(t)}, x) \; I(pa_i^{(t)}, \pi)}{\sum_t I(pa_i^{(t)}, \pi)}$$

\*    Learning from incomplete data

    examples $t = 1, 2, \ldots, T$
    Visible nodes   $V^{(t)}$
    hidden nodes   $H^{(t)}$

    Maximize likelihood with EM algorithm:

$$P(x_i = x \mid pa_i = \pi) \longleftarrow \frac{\sum_{t=1}^{T} P(x_i = x, \, pa_i = \pi \mid V^{(t)})}{\sum_{t=1}^{T} P(pa_i = \pi \mid V^{(t)})}$$
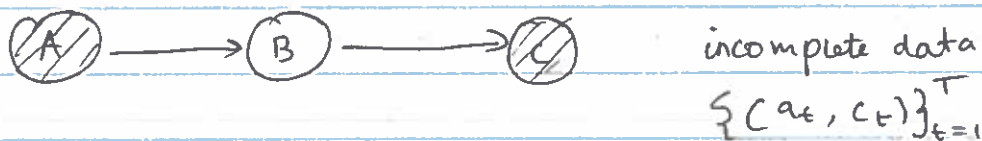
(new CPTs)

(computed from old CPTs)

Ex:      complete data $\{(a_t, b_t, c_t)\}_{t=1}^{T}$

$$P_{ML}(B = b \mid A = a) = \frac{count(A = a, \, B = b)}{count(A = a)}$$

$$= \frac{\sum_t I(a^{(t)}, a) \, I(b^{(t)}, b)}{\sum_t I(a^{(t)}, a)}$$

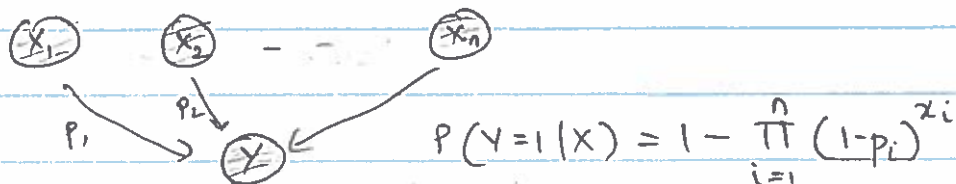   incomplete data $\{(a_t, c_t)\}_{t=1}^{T}$

$$P(B = b \mid A = a) \longleftarrow \frac{\sum_t I(a^{(t)}, a) \, \boxed{P(b \mid a_t \, c_t)}}{\sum_t I(a^{(t)}, a)}$$

(new CPTs)

computed via Bayes rule using current CPTs

\*    EM algorithm for noisy-OR model



$$P(Y = 1 \mid X) = 1 - \prod_{i=1}^{n} (1 - p_i)^{x_i}$$

Equivalent to:



$$P(z_i = 0 | x_i) = (1-p_i)^{x_i}$$

$$y = \text{logical-OR} (z_1, z_2 \ldots z_n)$$

$P(Y=1 | \bar{x})$ same form as noisy-OR

---

| Markov models of language |
| --- |

Let $w_l = l^{th}$ word in sentence. How to model $P(w_1 w_2 \ldots w_L)$? Shorthand $\vec{w} = (w_1, w_2 \ldots w_L)$

| Model | $P(\vec{w})$ | ML estimate | DAG |
| --- | --- | --- | --- |
| unigram | $\prod_l P_1(w_l)$ | $P_1(w) = \dfrac{\text{count}(w)}{\sum_{w'} \text{count}(w')}$ | $\textcircled{w}$ $\textcircled{w}$ $\cdots$ $\textcircled{w_L}$ |
| bigram | $P_1(w_1) \prod_{l=1} P(w_l | w_{l-1})$ | $P_2(w'|w) = \dfrac{\text{count}(w \to w')}{\text{count}(w)}$ | $\textcircled{w_1} \to \textcircled{w_2} \to \cdots \textcircled{w_L}$ |
| trigram | | | |

* Evaluating n-grams

- Train on corpus A, $P_1(\vec{w}) \leq P_2(\vec{w})$ on corpus A.
- Test on corpus B, $P_1(\vec{w}) \geq P_2(\vec{w})$ on corpus B

especially if $P_2(\vec{w}) = 0$ (unseen bigrams)

## Linear interpolation

\* Also known as mixture model

$$P_M(w_\ell \mid w_{\ell-1}) = \lambda P_1(w_\ell) + (1-\lambda) P_2(w_\ell \mid w_{\ell-1})$$

How to estimate $\lambda$ ?

\* Methodology

- Train $P_1, P_2$ on corpus A          $A =$ 'training set'
- fix $P_1$ and $P_2$              $B =$ 'test set'
- estimate $\lambda$ on corpus C        $C =$ 'development set'

Choose $\lambda$ to maximize likelihood $\prod_{\ell=1}^{L} P_M(w_\ell \mid w_{\ell-1})$ on corpus C.

Why not estimate $\lambda$ on other corpora ?
- Don't estimate $\lambda$ on B (test) — cheating
- Don't estimate $\lambda$ on A (train) — this would always yield $\lambda = 0$

\* Hidden variable model

$$P(w_\ell \mid w_{\ell-1}, z) = \begin{cases} P_1(w_\ell) & \text{if } z=1 \\ P_2(w_\ell \mid w_{\ell-1}) & \text{if } z=2 \end{cases}$$

$w_{\ell-1} \longrightarrow w_\ell$

$z = \{1,2\}$   $z$   $P(z=1) = \lambda$
$P(z=2) = 1-\lambda$

How to estimate $\lambda$ on corpus $C$?

$\{(w_{\ell-1}, w_\ell)\}_{\ell=1}^{L}$    incomplete data

In this model:

$$P(w_\ell | w_{\ell-1}) = \sum_{z=1}^{2} P(w_\ell, z | w_{\ell-1}) \quad \text{marginalization}$$

$$= \sum_{z} P(z | w_{\ell-1}) P(w_\ell | z, w_{\ell-1}) \quad \text{product rule}$$

$$= \sum_{z} P(z) P(w_\ell | z, w_{\ell-1}) \quad \text{independence}$$

$$= P(z=1) P(w_\ell | w_{\ell-1}, z=1) + P(z=2) P(w_\ell | w_{\ell-1}, z=2)$$

$$= \lambda P_1(w_\ell) + (1-\lambda) P_2(w' | w) \quad \begin{array}{l} \text{matches mixture} \\ \text{model} \end{array}$$

\*   <u>E-step</u>: compute posterior prob data

$$P(z | w_{\ell-1}, w_\ell) = \frac{P(w_\ell | z, w_{\ell-1}) P(z | w_{\ell-1})}{P(w_\ell | w_{\ell-1})} \quad \begin{array}{l} \text{Bayes,} \\ \text{marginal} \\ \text{ind.} \end{array}$$

$$P(z=1 | w_{\ell-1}, w_\ell) = \frac{\lambda P_1(w_\ell)}{\lambda P_1(w_\ell) + (1-\lambda) P_2(w_\ell | w_{\ell-1})}$$

$$P(z=2 | w_{\ell-1}, w_\ell) = 1 - P(z=1 | w_\ell, w_{\ell-1})$$

* M-step of EM algorithm

General rule:

$$P(\text{child}=c \mid pa = \pi) \leftarrow \frac{\sum_t P(\text{child}=c, pa_i = \pi \mid V^{(t)})}{\sum_t P(pa = \pi \mid V^{(t)})}$$

In this model:

$$P(z=1) \leftarrow \frac{\sum_{l=1}^{L} P(z=1 \mid w_{l-1}, w_l)}{L}$$

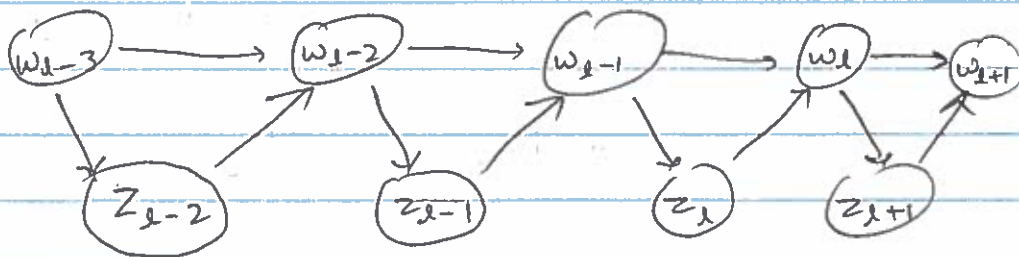$$\boxed{\lambda \leftarrow \frac{1}{L} \sum_l P(z=1 \mid w_{l-1}, w_l)}$$

* Iterate EM:

$$\lambda \leftarrow \frac{1}{L} \sum P(z=1 \mid w_l, w_{l-1})$$

guarantees improvement on corpus ( of log-likelihood

$$\mathcal{L}(\lambda) = \sum_l \log P_M(w_l \mid w_{l-1})$$

* In real-world application, mixing parameter would depend on previous word.



old $P(z_l = 1) = \lambda$

new $P(z_l = 1 \mid w_{l-1}) = \lambda_{w_{l-1}}$

* EM used to estimate as many params as words in vocab.

## Hidden Markov models  (HMMS)

**#** Variables

$S_t \in \{1, 2, \ldots n\}$    hidden state at time $t$

$O_t \in \{1, 2, \ldots m\}$    observation at time $t$

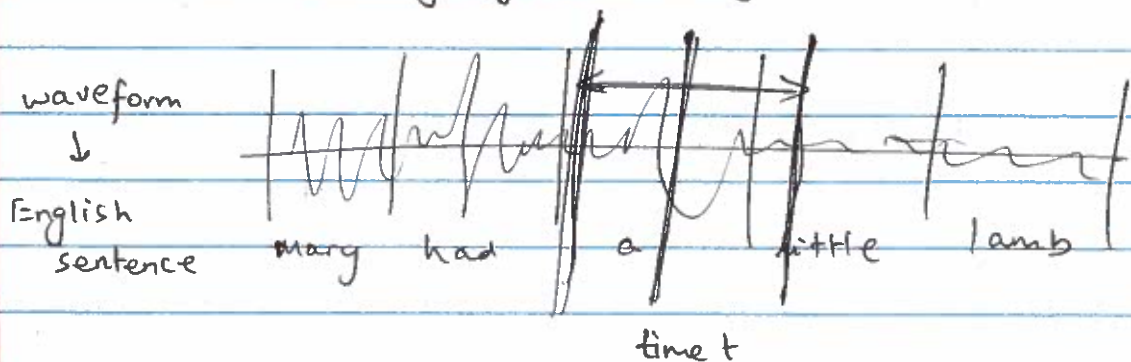(partial, noisy reflection of underlying hidden state)

Ex: house-training puppy

$S_t \in \{$ have to go, doesn't need to, went $\}$

$O_t \in \{$ wagging tail, barking, whimpering, running ... $\}$

Ex: speech recognition

$S_t$ = units of language (words, syllables, phonemes ...)

waveform
↓
English
sentence    Mary   had    a    little    lamb

time $t$

$O_t$ = acoustic measurements in sliding window centered at time $t$.

## Ex : robotics

$S_t$ = location, orientation, etc of robot at time $t$.

$O_t$ = sensor readings at time $t$

(camera, radar, heat map, ... )

### Markov assumptions

- finite context

$$P(S_t | S_1, S_2 \ldots S_{t-1}) = P(S_t | S_{t-1})$$

$$P(O_t | S_1, S_2 \ldots S_t, S_{t+1}, \ldots S_T) = P(O_t | S_t)$$

$$\underbrace{\hspace{5cm}}_{\text{all of time}}$$
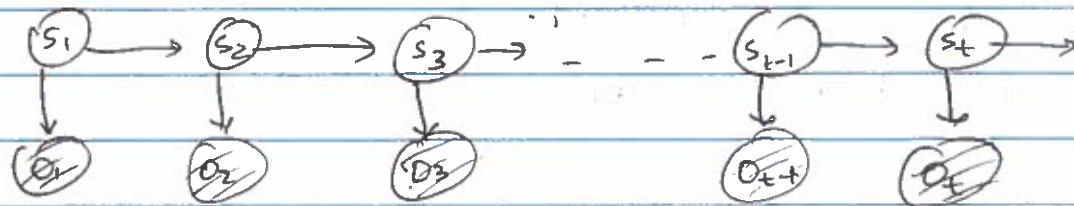
- shared CPTs

$$P(S_{t+1} = s' | S_t = s) = P(S_{t+k+1} = s' | S_{t+k} = s)$$

$$P(O_t = o | S_t = s) = O(_{t+k} = o | S_{t+k} = s)$$

### Belief network

$S_t \in \{1, 2, \ldots n\}$ hidden states

$O_t \in \{1, 2, \ldots m\}$ observations



Is it a polytree? Yes!

Joint distribution

$$P(\vec{s}, \vec{o}) = P(s_i) \left\{ \prod_{t=2}^{T} P(s_t \mid s_{t-1}) \right\} \left\{ \prod_{t=1}^{T} P(o_t \mid s_t) \right\}$$

$$\vec{s} = (s_1, s_2 \dots s_T)$$
$$\vec{o} = (o_1, o_2 \dots o_T)$$ shorthand

Parameters:

$\pi_i = P(s_1 = i)$   initial state distribution

$a_{ij} = P(s_{t+1} = j \mid s_t = i)$   transition matrix $(n \times n)$

$b_{ik} = P(o_t = k \mid s_t = i)$   emission matrix $(n \times m)$

Key questions

1) How to compute likelihood $P(o_1, o_2 \dots o_{T-1}, o_T)$?

2) How to compute most likely hidden state sequence?

$$\underset{\vec{s}}{\arg\max} \left\{ P(\underbrace{s_1, s_2 \dots s_T}_{n^T \text{ possible state sequences of length } T} \mid o_1, o_2 \dots o_T) \right\} = (s_1^*, s_2^* \dots s_T^*)$$

$$\vec{s} = (s_1, s_2 \dots s_T)$$

3) How to update beliefs for real-time monitoring?
How to compute $P(s_t = i \mid o_1, o_2 \dots o_{t-1}, o_t)$?

4) How to learn an HMM from data? How to estimate $\{\pi_i, a_{ij}, b_{ik}\}$ to maximize $P(o_1, o_2 \dots o_T)$?

(Use EM algorithm)

Questions 1-3 are inference for fixed HMM with given $\{\pi_i, a_{ij}, b_{ik}\}$

Question 4 - learning

___

i) How to compute likelihood?

marginalization

$$P(O_1, O_2 \dots O_T) = \sum_{\vec{S}} P(S_1, S_2 \dots S_T, O_1, O_2 \dots O_T)$$

$\rightarrow$ sum of $n^T$ sequences of hidden states

$$= \sum_{\vec{S}} P(S_1) \left\{ \prod_{t=2}^{T} P(S_t | S_{t-1}) \right\} \left\{ \prod_{t=1}^{T} P(O_t | S_t) \right\}$$

\* Efficient recursion

$$P\left(O_1, O_2 \dots O_{t-1}, O_t, \cancel{\Sigma} O_{t+1}, S_{t+1} = j\right)$$

$$= \sum_{i=1}^{n} P(O_1, O_2 \dots O_t, O_{t+1}, S_t = i, S_{t+1} = j) \quad \text{(marginalization)}$$

$$= \sum_{i=1}^{n} P(\underbrace{O_1, O_2 \dots O_t, S_t = i}_{\substack{\text{up to} \\ \text{time } t}}) \, P(S_{t+1} = j | S_t = i, O_1 \dots O_t) \times \overset{\text{(product}}{\text{rule)}}$$
$$P(O_{t+1} | S_{t+1} = j, S_t = i, O_1 \dots O_t)$$

$$\overset{(\underline{+})}{=} \sum_{i=1}^{n} \underbrace{P(O_1 \dots O_t, S_t = i)}_{\substack{\text{recursive} \\ \text{instance}}} \underbrace{P(S_{t+1} = j | S_t = i)}_{a_{ij}} \underbrace{P(O_{t+1} | S_{t+1} = j)}_{b_{j, O_{t+1}}}$$

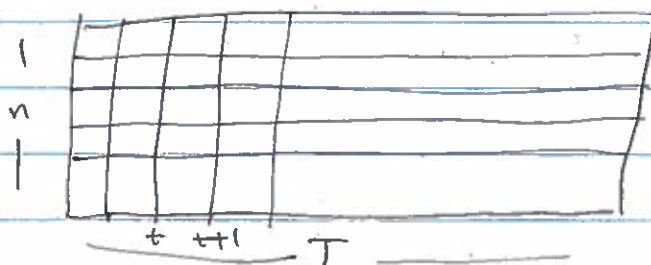$$\underbrace{\hspace{6cm}}_{\text{CPTs}}$$

\* shorthand

$$\alpha_{it} = P(o_1 \dots o_t, s_t = i)$$

$\quad\hookrightarrow i = 1 \dots n$ # hidden states

$\quad\quad t = 1 \dots T$ sequence length.



\* Forward algorithm in HMMs

− base case $(t=1)$ first column of matrix

$$\alpha_{i1} = P(o_1, s_1 = i) \quad \text{by definition}$$
$$= P(s_1 = i) \, P(o_1 | s_1 = i)$$
$$= \pi_i \, \underline{b_i(o_1)}$$

$\quad\quad\quad\quad\quad$ emission matrix
$\quad\quad\quad\quad\quad$ element $b_{io_1}$

− recursion step from time $t$ to $t+1$

$$\alpha_{j, t+1} = \sum_{i=1}^{n} \alpha_{it} \, a_{ij} \, b_j(o_{t+1})$$

\* Back to likelihood : $\hspace{3cm}$ marginalization

$$P(o_1, o_2 \dots o_T) = \sum_{i=1}^{n} P(o_1, o_2 \dots o_T, s_T = i)$$

$$= \sum_{i=1}^{\hat{n}} \alpha_{iT}$$

\*    Scales as $O(Tn^2)$

     - linear, not exponential, in sequence length.

     - quadratic in # hidden states $n$

\*    Warning: naive calculation with underflow for long sequences $T \gg 1$ because $P(O_1 O_2 .. O_T) \ll 1$

     (fix by rescaling at each iteration computing log-likelihood)

2)   How to compute most likely state sequence?

$$\vec{s}^* = \{ s_1^*, s_2^* .... s_T^* \}$$

$$= \underset{\vec{s}}{argmax} \left[ P(s_1 \, s_2 ... s_T \mid O_1 \, O_2 - - O_T) \right]$$

$$= \underset{\vec{s}}{argmax} \left[ \frac{P(s_1, s_2 .. s_T, O_1, ... O_T)}{P(O_1, O_2 - - O_T)} \right] \quad \text{product rule}$$

$$\hookrightarrow \text{ind of } \vec{s}$$

$$= \underset{\vec{s}}{argmax} \left[ P(s_1 \, s_2 - - s_T, O_1, - - O_T) \right]$$
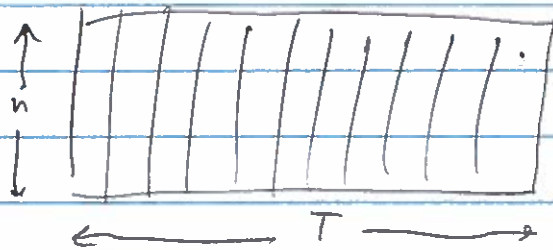
         How to compute $s^*$?

Define :

$$\ell^*_{it} = \max_{s_1, s_2 \cdots s_{t-1}} \left\{ \log P \left( s_1, s_2 \cdots s_{t-1}, s_t = i, 0_1 \cdots 0_t \right) \right\}$$

$i = 1 \cdots n$  # hidden states

$t = 1 \cdots T$ sequence length



$\longleftarrow \quad T \quad \longrightarrow$

$= \log\text{-probability of most}$
likely sub-sequence of
hidden states $s_1, s_2 \cdots s_t$
that ends in state $s_t = i$
and explains observations
$0_1, 0_2 \cdots 0_t$