

LECTURE 18:

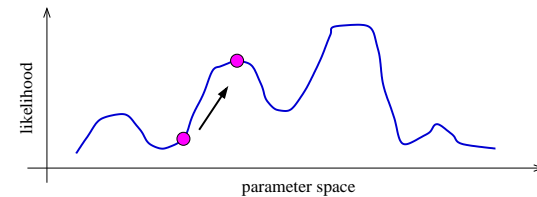
HIDDEN MARKOV MODEL LEARNING

Sam Roweis

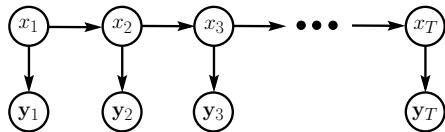
March 10, 2004

BAUM-WELCH ALGORITHM: EM TRAINING

1. Intuition: if only we *knew* the true state path then ML parameter estimation would be trivial (MM1 on x , conditional on y).
2. But: can *estimate* state path using inference recursions.
3. *Baum-Welch algorithm* (special case of EM): estimate the states, then compute params, then re-estimate states, and so on ...
4. This works and we can *prove* that it always improves likelihood.
5. However: finding the ML parameters is NP complete, so initial conditions matter a lot and convergence is hard to tell.



REMINDER: HMM GRAPHICAL MODEL



- Hidden states $\{x_t\}$, outputs $\{y_t\}$

Joint probability factorizes:

$$\begin{aligned} P(\{x\}, \{y\}) &= \prod_{t=1}^T P(x_t|x_{t-1})P(y_t|x_t) \\ &= \pi_{x_1} \prod_{t=1}^{T-1} S_{x_t, x_{t+1}} \prod_{t=1}^T A_{x_t}(y_t) \end{aligned}$$

- We saw efficient recursions for computing

$$L = P(\{y\}) = \sum_{\{x\}} P(\{x\}, \{y\}) \text{ and } \gamma_i(t) = P(x_t = i|\{y\}).$$

PARAMETER ESTIMATION USING EM

- S_{ij} are transition probs; state j has output distribution $A_j(y)$

$$P(x_{t+1} = j|x_t = i) = S_{ij} \quad P(x_1 = j) = \pi_j$$

$$P(y_t = y|x_t = j) = A_j(y)$$

- Complete log likelihood:

$$\begin{aligned} \log p(x, y) &= \log \left\{ \pi_{x_1} \prod_{t=1}^{T-1} S_{x_t, x_{t+1}} \prod_{t=1}^T A_{x_t}(y_t) \right\} \\ &= \log \left\{ \prod_i \pi_i^{[x_1^i]} \prod_{t=1}^{T-1} \prod_{ij} S_{ij}^{[x_t^i, x_{t+1}^j]} \prod_{t=1}^T \prod_k A_k(y_t)^{[x_t^k]} \right\} \\ &= \sum_i [x_1^i] \log \pi_i + \sum_{t=1}^{T-1} \sum_{ij} [x_t^i, x_{t+1}^j] \log S_{ij} + \sum_{t=1}^T \sum_k [x_t^k] \log A_k(y_t) \end{aligned}$$

where the indicator $[x_t^i] = 1$ if $x_t = i$ and 0 otherwise

- For EM, we need to compute the *expected complete log likelihood*.

STATE EXPECTATIONS REQUIRED FROM THE E-STEP

- The expected complete log likelihood requires
 $\gamma_i(t) = \langle [x_t^i] \rangle$ and $\xi_{ij}(t) = \langle [x_t^i, x_{t+1}^j] \rangle$
- So in the E-step we need to compute both
 $\gamma_i(t) = p(x_t = i | \{\mathbf{y}\})$ and $\xi_{ij}(t) = p(x_t = i, x_{t+1} = j | \{\mathbf{y}\})$.
- We already know how to compute $\gamma_i(t)$ using α and β recursions.
 We can compute $\xi_{ij}(t)$ the same way (recall BP):

$$\begin{aligned} \xi_{ij}(t) &= p(x_{it}, x_{jt+1} | \{\mathbf{y}\}) = p(x_{it} | \{\mathbf{y}\}) p(x_{jt+1} | x_{it}, \{\mathbf{y}\}) \\ &= p(x_{it}, y_1^t | y_{t+1}^T) p(x_{jt+1} | x_{it}, y_{t+1}^T) / p(y_1^t | y_{t+1}^T) \\ &= \frac{p(x_{it}, y_1^t) p(y_{t+1}^T | x_{it}, y_1^t) p(y_{t+1}^T | x_{jt+1}, x_{it}) p(x_{jt+1} | x_{it})}{p(y_1^t | y_{t+1}^T) p(y_{t+1}^T | x_i = t)} \\ &= \frac{p(x_{it}, y_1^t) p(y_{t+1}^T | x_{it}) p(y_{t+1} | x_{jt+1}) p(y_{t+2}^T | x_{jt+1}) p(x_{jt+1} | x_{it})}{p(y_1^T) p(y_{t+1}^T | x_i = t)} \\ &= \alpha_i(t) A_j(y_{t+1}) S_{ij} \beta_j(t+1) / L \end{aligned}$$

HMM PRACTICALITIES

- Multiple observation sequences: can be dealt with by averaging numerators and averaging denominators in the ratios given above.
- Initialization: mixtures of Naive Bayes or mixtures of Gaussians
- Numerical scaling: the probability values that the bugs carry get tiny for big times and so can easily underflow. Good rescaling trick:

$$\rho_t = P(\mathbf{y}_t | \mathbf{y}_1^{t-1}) \quad \alpha(t) = \bar{\alpha}(t) \prod_{t'=1}^t \rho_{t'}$$

or represent all probabilities as logs and use logsum

M-STEP: NEW PARAMETERS ARE JUST RATIOS OF FREQUENCY COUNTS

- Initial state distribution: expected #times in state i at time 1:

$$\hat{\pi}_i = \gamma_i(1)$$

- Expected #transitions from state i to j which begin at time t :

$$\xi_{ij}(t) = \alpha_i(t) S_{ij} A_j(\mathbf{y}_{t+1}) \beta_j(t+1) / L$$

so the estimated transition probabilities are:

$$\hat{S}_{ij} = \frac{\sum_{t=1}^{T-1} \xi_{ij}(t)}{\sum_{t=1}^{T-1} \gamma_i(t)}$$

- The output distributions are the expected number of times we observe a particular symbol in a particular state:

$$\hat{A}_j(y_0) = \frac{\sum_{t | \mathbf{y}_t = y_0} \gamma_j(t)}{\sum_{t=1}^T \gamma_j(t)}$$

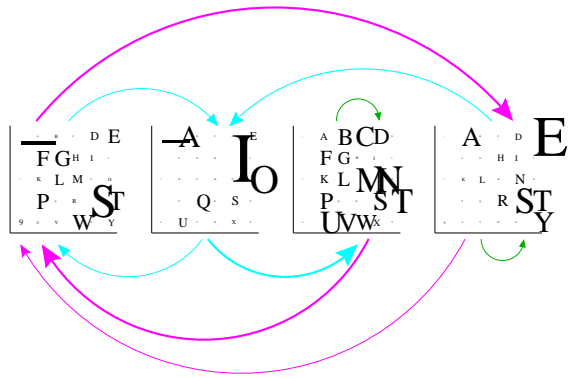
M-STEP FOR PROFILE HMMs

- The emission probabilities $A_j()$ for match and insert states and the initial state distribution π (for m_1, i_1, d_1) are updated exactly as in the regular M-step.
- The expected #transitions from state i to j which begin at time t are different when j is a delete state:

$$\xi_{ij}(t) = \alpha_i(t) S_{ij} \beta_j(t) / L$$
- Given this change, the updates to the transition parameters is the same as in the normal M-step.

SYMBOL HMM EXAMPLE

- Character sequences (discrete outputs)



MIXTURE HMM EXAMPLE

- Geysers data (continuous outputs)

