

SCIA 2003 TUTORIAL:
HIDDEN MARKOV MODELS

Sam Roweis, University of Toronto

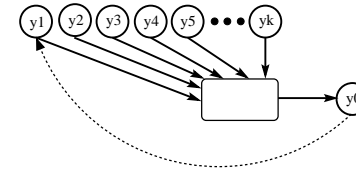
June 29, 2003

MARKOV MODELS

- Use past as state. Next output depends on previous output(s):

$$\mathbf{y}_t = f[\mathbf{y}_{t-1}, \mathbf{y}_{t-2}, \dots]$$

order is number of previous outputs



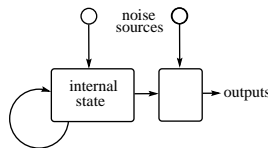
- Add noise to make the system stochastic:

$$p(\mathbf{y}_t | \mathbf{y}_{t-1}, \mathbf{y}_{t-2}, \dots, \mathbf{y}_{t-k})$$

- Markov models have two problems:
 - need big order to remember past “events”
 - output noise is confounded with state noise

PROBABILISTIC GENERATIVE MODELS FOR TIME SERIES

- Stochastic models for time-series: y_1, y_2, \dots, y_T
To get interesting variability need *noise*.
To get correlations across time, need some system *state*.



- Time, States, Outputs: can be discrete or continuous
- Today: discrete state, discrete time
similar to finite state automata; Moore/Mealy machines
- Main idea of generative models: invent a model for generating data, and set its parameters so it is as likely as possible to generate the data you actually observed.

LEARNING MARKOV MODELS

- The ML parameter estimates for a simple Markov model are easy:

$$p(\mathbf{y}_1, \dots, \mathbf{y}_T) = p(\mathbf{y}_1 \dots \mathbf{y}_k) \prod_{t=k+1}^T p(\mathbf{y}_t | \mathbf{y}_{t-1}, \dots, \mathbf{y}_{t-k})$$

$$\log p(\{\mathbf{y}\}) = \log p(\mathbf{y}_1 \dots \mathbf{y}_k) + \sum_{t=k+1}^T \log p(\mathbf{y}_t | \mathbf{y}_{t-1}, \mathbf{y}_{t-2}, \dots, \mathbf{y}_{t-k})$$

- Each window of $k + 1$ outputs is a training case for the model $p(\mathbf{y}_t | \mathbf{y}_{t-1}, \mathbf{y}_{t-2}, \dots, \mathbf{y}_{t-k})$.
- Example: for discrete outputs (symbols) and a 2nd-order markov model we can use the multinomial model:

$$p(y_t = m | y_{t-1} = a, y_{t-2} = b) = \alpha_{mab}$$

The maximum likelihood values for α are:

$$\alpha_{mab}^* = \frac{\text{num}[t \text{ s.t. } y_t = m, y_{t-1} = a, y_{t-2} = b]}{\text{num}[t \text{ s.t. } y_{t-1} = a, y_{t-2} = b]}$$

HIDDEN MARKOV MODELS (HMMS)

Add a latent (hidden) variable x_t to improve the model.

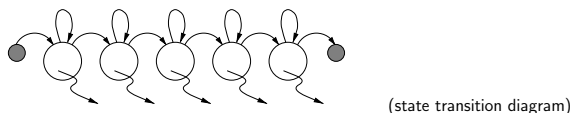
- HMM \equiv "probabilistic function of a Markov chain":

1. 1st-order Markov chain generates hidden state sequence (path):

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

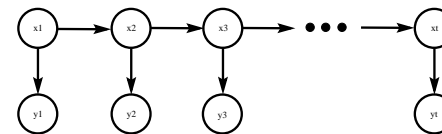
2. A set of output probability distributions $\mathbf{A}_j(\cdot)$ (one per state) converts state path into sequence of observable symbols/vectors

$$p(\mathbf{y}_t = y | x_t = j) = \mathbf{A}_j(\mathbf{y})$$



- Even though hidden state seq. is 1st-order Markov, the output process is not Markov of *any* order [ex. 1111121111311121111131...]

HMM MODEL EQUATIONS



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

Joint probability factorizes:

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

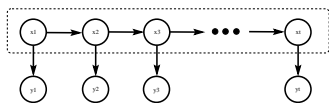
- NB: Data are *not* i.i.d. Everything is coupled across time.
- Three problems: computing probabilities of observed sequences, inference of hidden state sequences, learning of parameters.

LINKS TO OTHER MODELS

- You can think of an HMM as:

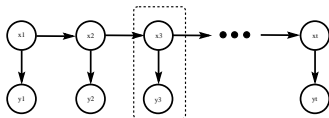
A Markov chain with stochastic measurements.

(graphical models)



or

A mixture model with states coupled across time.



- The future is independent of the past given the present. However, conditioning on all the observations couples hidden states.

PROBABILITY OF AN OBSERVED SEQUENCE

- To evaluate the probability $p(\{y\})$, we want:

$$p(\{y\}) = \sum_{\{x\}} p(\{x\}, \{y\})$$

$$p(\text{observed sequence}) = \sum_{\text{all paths}} p(\text{observed outputs, state path})$$

- Looks hard! (#paths = #states^T).

But joint probability factorizes:

$$\begin{aligned} p(\{y\}) &= \sum_{x_1} \sum_{x_2} \cdots \sum_{x_T} \prod_{t=1}^T p(x_t | x_{t-1}) p(y_t | x_t) \\ &= \sum_{x_1} p(x_1) p(y_1 | x_1) \sum_{x_2} p(x_2 | x_1) p(y_2 | x_2) \cdots \\ &\quad \sum_{x_T} p(x_T | x_{T-1}) p(y_T | x_T) \end{aligned}$$

- By moving the summations inside, we can save a lot of work.

IT'S AS EASY AS ABC...

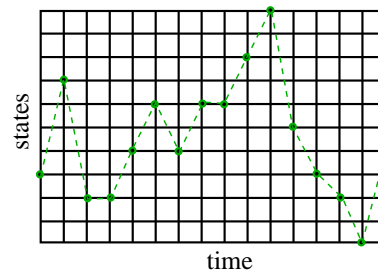
If you understand this:

$$ab + ac = a(b + c)$$

then you understand the main trick of HMMs.

BUGS ON A GRID

- Naive algorithm:
 1. start bug in each state at $t=1$ holding value 0
 2. move each bug forward in time by making copies of it and incrementing the value of each copy by the probability of the transition and output emission
 3. go to 2 until all bugs have reached time T
 4. sum up values on all bugs



THE SUM-PRODUCT RECURSION

- We want to compute:

$$L = p(\{\mathbf{y}\}) = \sum_{\{x\}} p(\{x\}, \{\mathbf{y}\})$$

- There exists a clever “forward recursion” to compute this huge sum very efficiently. Define $\alpha_j(t)$:

$$\alpha_j(t) = p(\mathbf{y}_1^t, x_t = j)$$

$$\alpha_j(1) = \pi_j \mathbf{A}_j(\mathbf{y}_1) \quad \text{now induction to the rescue...}$$

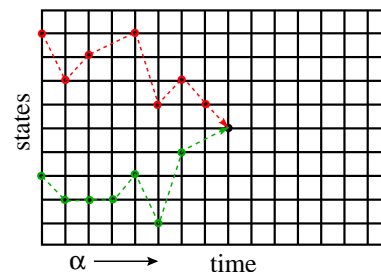
$$\alpha_k(t+1) = \left\{ \sum_j \alpha_j(t) S_{jk} \right\} A_k(\mathbf{y}_{t+1})$$

- Notation: $x_a^b \equiv \{x_a, \dots, x_b\}$; $\mathbf{y}_a^b \equiv \{\mathbf{y}_a, \dots, \mathbf{y}_b\}$
- This enables us to easily (cheaply) compute the desired likelihood L since we know we must end in some possible state:

$$L = \sum_k \alpha_k(T)$$

BUGS ON A GRID - TRICK

- Clever recursion:
 - adds a step between 2 and 3 above which says: at each node, replace all the bugs with a single bug carrying the sum of their values



- This trick is called *dynamic programming*, and can be used whenever we have a summation, search, or maximization problem that can be set up as a grid in this way. The axes of the grid don't have to be “time” and “states”.

INFERENCE OF HIDDEN STATES

- What if we want to estimate the hidden states given observations? To start with, let us estimate a single hidden state:

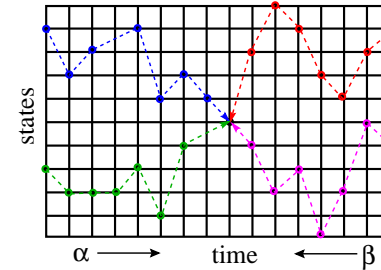
$$\begin{aligned}
 p(x_t|\{\mathbf{y}\}) &= \gamma(x_t) = \frac{p(\{\mathbf{y}\}|x_t)p(x_t)}{p(\{\mathbf{y}\})} \\
 &= \frac{p(\mathbf{y}_1^t|x_t)p(\mathbf{y}_{t+1}^T|x_t)p(x_t)}{p(\mathbf{y}_1^T)} \\
 &= \frac{p(\mathbf{y}_1^t, x_t)p(\mathbf{y}_{t+1}^T|x_t)}{p(\mathbf{y}_1^T)} \\
 p(x_t|\{\mathbf{y}\}) &= \gamma(x_t) = \frac{\alpha(x_t)\beta(x_t)}{p(\mathbf{y}_1^T)}
 \end{aligned}$$

where

$$\begin{aligned}
 \alpha_j(t) &= p(\mathbf{y}_1^t, x_t = j) \\
 \beta_j(t) &= p(\mathbf{y}_{t+1}^T | x_t = j) \\
 \gamma_i(t) &= p(x_t = i | \mathbf{y}_1^T)
 \end{aligned}$$

FORWARD-BACKWARD ALGORITHM

- $\alpha_i(t)$ gives total *inflow* of prob. to node (t, i)
 $\beta_i(t)$ gives total *outflow* of prob.



- Bugs again: we just let the bugs run forward from time 0 to t and backward from time T to t .
- In fact, we can just do one forward pass to compute all the $\alpha_i(t)$ and one backward pass to compute all the $\beta_i(t)$ and then compute any $\gamma_i(t)$ we want. Total cost is $O(K^2T)$.

FORWARD-BACKWARD ALGORITHM

- We compute these quantities efficiently using another recursion. Use total prob. of all paths going through state i at time t to compute the *conditional* prob. of being in state i at time t :

$$\begin{aligned}
 \gamma_i(t) &= p(x_t = i | \mathbf{y}_1^T) \\
 &= \alpha_i(t)\beta_i(t)/L
 \end{aligned}$$

where we defined:

$$\beta_j(t) = p(\mathbf{y}_{t+1}^T | x_t = j)$$

- There is also a simple recursion for $\beta_j(t)$:

$$\begin{aligned}
 \beta_j(t) &= \sum_i S_{ji}\beta_i(t+1)A_i(\mathbf{y}_{t+1}) \\
 \beta_j(T) &= 1
 \end{aligned}$$

- $\alpha_i(t)$ gives total *inflow* of prob. to node (t, i)
 $\beta_i(t)$ gives total *outflow* of prob.

LIKELIHOOD FROM FORWARD-BACKWARD ALGORITHM

- Since $\sum_{x_t} \gamma(x_t) = 1$, we can compute the likelihood at *any* time using the results of the $\alpha - \beta$ recursions:

$$L = p(\{\mathbf{y}\}) = \sum_{x_t} \alpha(x_t)\beta(x_t)$$

- In the forward calculation we proposed originally, we did this at the final timestep $t = T$:

$$L = \sum_{x_T} \alpha(x_T)$$

because $\beta_T = 1$.

- This is a good way to check your code!

VITERBI DECODING: MAX-PRODUCT

- The numbers $\gamma_j(t)$ above gave the probability distribution over all states at any time.
- By choosing the state $\gamma_*(t)$ with the largest probability at each time, we can make an “best” state path. This is the path with the *maximum expected number of correct states*.
- But it *is not* the single path with the highest likelihood of generating the data. In fact it may be a path of probability zero!
- To find the single best path, we do *Viterbi decoding* which is just Bellman's dynamic programming algorithm applied to this problem.
- The recursions look the same, except with \max instead of \sum .
- Bugs once more: same trick except at each step kill all bugs but the one with the highest value at the node.

PARAMETER ESTIMATION

- Complete log likelihood:

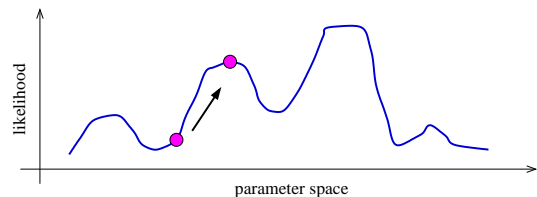
$$\begin{aligned} \log p(\{x\}, \{y\}) &= \log \left\{ \pi_{x_1} \prod_{t=1}^{T-1} S_{x_t, x_{t+1}} \prod_{t=1}^T A_{x_t}(\mathbf{y}_t) \right\} \\ &= \log \left\{ \prod_i \pi_i^{[x_1^i]} \prod_{t=1}^{T-1} \prod_j S_{ij}^{[x_t^i, x_{t+1}^j]} \prod_{t=1}^T \prod_k A_k(\mathbf{y}_t)^{[x_t^k]} \right\}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

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

- EM maximizes *expected value* of $\log p(\{x\}, \{y\})$ under $p(\{x\}|\{y\})$
So the statistics we need from the inference (E-step) are:
 $p(x_t|\{y\})$ and $p(x_t, x_{t+1}|\{y\})$.
- We saw how to get single time marginals $p(x_t|\{y\})$, but what about two-frame estimates $p(x_t, x_{t+1}|\{y\})$?

BAUM-WELCH TRAINING: EM ALGORITHM

1. How to find the parameters? Intuition: if only we *knew* the true state path then ML parameter estimation would be trivial.
2. But: can *estimate* state path using the DP trick.
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 hard, so initial conditions matter a lot and convergence is hard to tell.



TWO-FRAME INFERENCE

- Need the cross-time statistics for adjacent time steps:

$$\xi_{ij} = p(x_t = i, x_{t+1} = j | \{y\})$$

- This can be done by rewriting:

$$\begin{aligned} p(x_t, x_{t+1} | \{y\}) &= p(x_t, x_{t+1}, \{y\}) / p(\{y\}) \\ &= p(x_t, \mathbf{y}_1^t) p(x_{t+1}, \mathbf{y}_{t+1}^T | x_t, \mathbf{y}_1^t) / L \\ &= p(x_t, \mathbf{y}_1^t) p(x_{t+1} | x_t) p(\mathbf{y}_{t+1} | x_{t+1}) p(\mathbf{y}_{t+2}^T | x_{t+1}) / L \\ &= \alpha_i(t) S_{ij} \mathbf{A}_j(\mathbf{y}_{t+1}) \beta_j(t+1) / L \\ &= \xi_{ij} \end{aligned}$$

- This is the expected number of transitions from state i to state j that begin at time t , given the observations.
- It can be computed with the same α and β recursions.

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 :

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

so the estimated transition probabilities are:

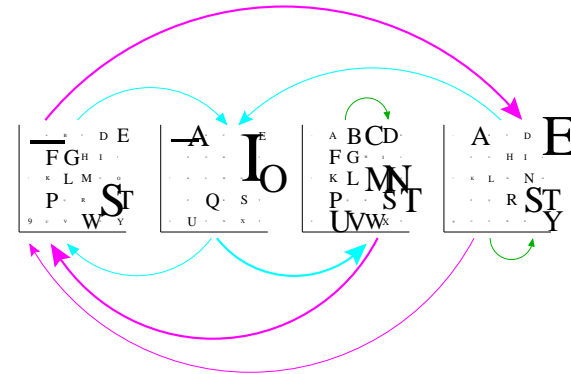
$$\hat{S}_{ij} = \frac{\sum_{t=1}^{T-1} \xi_{ij}(t)}{\sum_{t=1}^{T-1} \gamma_i(t)}$$

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

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

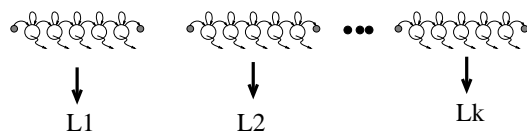
SYMBOL HMM EXAMPLE

- Character sequences (discrete outputs)



USING HMMs FOR RECOGNITION

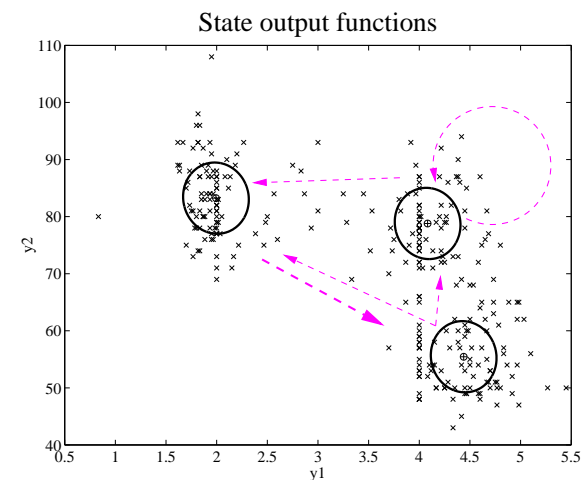
- Use many HMMs for recognition by:
 - training one HMM for each class (requires *labeled* training data)
 - evaluating probability of an unknown sequence under each HMM
 - classifying unknown sequence: HMM with highest likelihood



- This requires the solution of two problems:
 - Given model, evaluate prob. of a sequence.
(We can do this exactly & efficiently.)
 - Give some training sequences, estimate model parameters.
(We can find the local maximum of parameter space nearest our starting point using Baum-Welch (EM).)

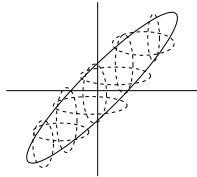
MIXTURE HMM EXAMPLE

- Geyser data (continuous outputs)



REGULARIZING HMMs

- Two problems:
 - for high dimensional outputs, lots of parameters in each $\mathbf{A}_j(\mathbf{y})$
 - with many states, transition matrix has many² elements
- First problem: full covariance matrices in high dimensions or discrete symbol models with many symbols have *lots* of parameters. To estimate these accurately requires a lot of training data. Instead, we often use mixtures of diagonal covariance Gaussians.



- For discrete data, use mixtures of base rates and/or smoothing.
- We can also tie parameters across states (shrinkage).

MORE ADVANCED TOPICS

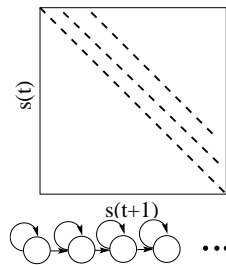
- Multiple observation sequences: can be dealt with by averaging numerators and averaging denominators in the ratios given above.
- Generation of new sequences. Just roll the dice!
- Sampling a single state sequence from the posterior $p(\{x\}|\{y\})$. Harder...but possible. (can you think of how?)
- Initialization: mixtures of base rates or mixtures of Gaussians. Can also use a trick of building a suffix tree to efficiently count all subsequences up to some length and using these counts to initialize.
- Null outputs: it is possible to have states which (sometimes or always) output nothing. This often makes the representation sparser (e.g. profile HMMs).
- There is also a modified Baum-Welch training based on the Viterbi decode. Like K-means instead of mixtures of Gaussians.

REGULARIZING TRANSITION MATRICES

- One way to regularize large transition matrices is to *constrain* them to be relatively *sparse*: instead of being allowed to transition to *any* other state, each state has only a few possible successor states.
- For example if each state has at most p possible next states then the cost of inference is $O(pKT)$ and the number of parameters is $O(pK + KM)$ which are both *linear* in the number of states.

An extremely effective way to constrain the transitions is to *order* the states in the HMM and allow transitions only to *states that come*

- *later in the ordering*. Such models are known as “linear HMMs”, “chain HMMs” or “left-to-right HMMs”. Transition matrix is upper-diagonal (usually only has a few bands).



BE CAREFUL: LOGSUM

- Often you can easily compute $b_k = \log p(y|x = k, \mathbf{A}_k)$, but it will be very negative, say -10^6 or smaller.
- Now, to compute $\ell = \log p(y|\mathbf{A})$ you need to compute $\log \sum_k e^{b_k}$.
- Careful! Do not compute this by doing $\log(\text{sum}(\exp(b)))$. You will get underflow and an incorrect answer.
- Instead do this:
 - Add a constant exponent B to all the values b_k such that the largest value comes close to the maximum exponent allowed by machine precision: $B = \text{MAXEXPONENT} - \log(K) - \max(b)$.
 - Compute $\log(\text{sum}(\exp(b+B))) - B$.
- Example: if $\log p(y|x = 1) = -120$ and $\log p(y|x = 2) = -120$, what is $\log p(y) = \log [p(y|x = 1) + p(y|x = 2)]$?
Answer: $\log[2e^{-120}] = -120 + \log 2$.

HMM PRACTICALITIES

- If you just implement things as I have described them, *they will not work at all*. Why? Remember logsum...
- Numerical scaling: the probability values that the bugs carry get tiny for big times and so can easily underflow. Good rescaling trick:

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

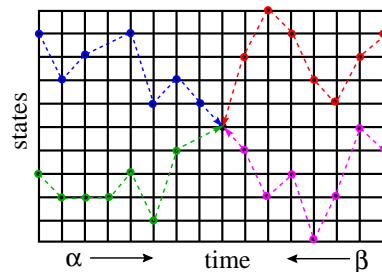
(note: some errors in early editions of Rabiner & Juang)

APPLICATIONS OF HMMs

- Speech recognition.
- Language modeling.
- Information retrieval.
- Motion video analysis/tracking.
- Protein sequence and genetic sequence alignment and analysis.
- Financial time series prediction.
- Modeling growth/diffusion of cells in biological systems.
- ...

COMPUTATIONAL COSTS IN HMMs

- The number of parameters in the model was $O(K^2 + KM)$ for M output symbols or dimensions.
- Recall the forward-backward algorithm for inference of state probabilities $p(x_t | \{\mathbf{y}\})$.
- The storage cost of this procedure was $O(KT + K^2)$ for K states and a sequence of length T .
- The time complexity was $O(K^2T)$.



LINEAR DYNAMICAL SYSTEMS (STATE SPACE MODELS)

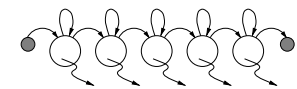
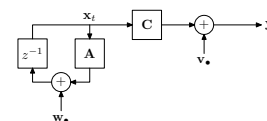
- LDS is a Gauss-Markov continuous state process:

$$x_{t+1} = \mathbf{A}x_t + w_t$$

observed through the "lens" of a noisy linear embedding:

$$\mathbf{y}_t = \mathbf{C}x_t + v_t$$

- Noises w_t and v_t are temporally white and uncorrelated
- Exactly the *continuous state analogue of Hidden Markov Models*.

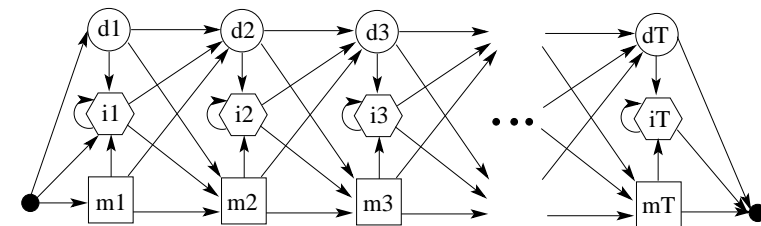


- forward algorithm \Leftrightarrow Discrete Kalman Filter
- forward-backward \Leftrightarrow Kalman Smoothing
- Viterbi decoding \Leftrightarrow no equivalent

DISCRETE SEQUENCES IN COMPUTATIONAL BIOLOGY

- There has recently been a great interest in applying probabilistic models to analyzing discrete sequence data in molecular and computational biology.
- There are two major sources of such data:
 - amino acid sequences for protein analysis
 - base-pair sequences for genetic analysis
- The sequences are sometimes annotated by other labels, e.g. species, mutation/disease type, gender, race, etc.
- Lots of interesting applications:
 - whole genome shotgun sequence fragment assembly
 - multiple alignment of conserved sequences
 - splice site detection
 - inferring phylogenetic trees

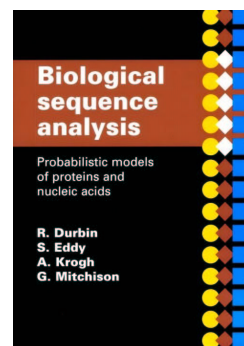
PROFILE (STRING-EDIT) HMMs



- A “profile HMM” or “string-edit” HMM is used for probabilistically matching an observed input string to a stored template pattern with possible insertions and deletions.
- Three kinds of states:
 - m_j – use position j in the template to match an observed symbol
 - i_j – insert extra symbol(s) observations after template position j
 - d_j – delete (skip) template position j

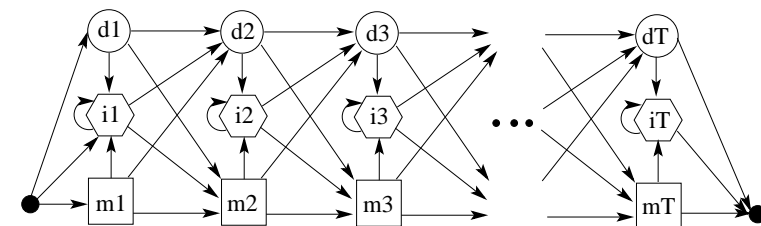
MAIN TOOL: HIDDEN MARKOV MODELS

- HMMs and related models (e.g. profile HMMs) have been the major tool used in biological sequence analysis and alignment.
- The basic dynamic programming algorithms can be improved in special cases to make them more efficient in time or memory.



See the excellent book by Durbin, Eddy, Krogh, Mitchison for lots of practical details on applications and implementations.

PROFILE HMMs HAVE LINEAR COSTS



- number of states = $3(\text{length_template})$
- Only insert and match states can generate output symbols.
- Once you visit or skip a match state you can never return to it.
- At most 3 destination states from any state, so S_{ij} very sparse.
- Storage/Time cost *linear* in #states, not quadratic.
- State variables and observations no longer in sync.
(e.g. $y_1:m_1 ; d_2 ; y_2:i_2 ; y_3:i_2 ; y_4:m_3 ; \dots$)

