

CSCE
471/871

Lecture 4:
Profile Hidden
Markov
Models

Stephen Scott

Organization

Building
Models

Searching

CSCE 471/871 Lecture 4: Profile Hidden Markov Models

Stephen Scott

sscott@cse.unl.edu

CSCE
471/871

Lecture 4:
Profile Hidden
Markov
Models

Stephen Scott

Organization

Building
Models

Searching

- Designed to model (profile) a multiple alignment of a protein family (e.g., Fig. 5.1)
- Gives a probabilistic model of the proteins in the family
- Useful for searching databases for more homologues and for aligning strings to the family

CSCE
471/871

Lecture 4:
Profile Hidden
Markov
Models

Stephen Scott

Organization

Building
Models

Searching

- Organization of a profile HMM
 - Ungapped regions
 - Insert and delete states
 - Non-global alignments
- Building a model
 - Determining states: match, insert, delete
 - Estimating probabilities
 - Pseudocounts
- Searching and aligning with HMMs
 - Viterbi
 - Forward

Organization of a Profile HMM

Match States

CSCE
471/871

Lecture 4:
Profile Hidden
Markov
Models

Stephen Scott

Organization

Match States

Insertion States

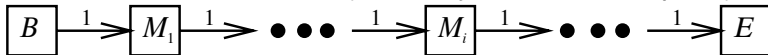
Deletion States

General Structure

Building
Models

Searching

Start with a trivial HMM M (not really hidden at this point)



Each match state has its own set of emission probabilities, so we can compute probability of a new sequence x being part of this family:

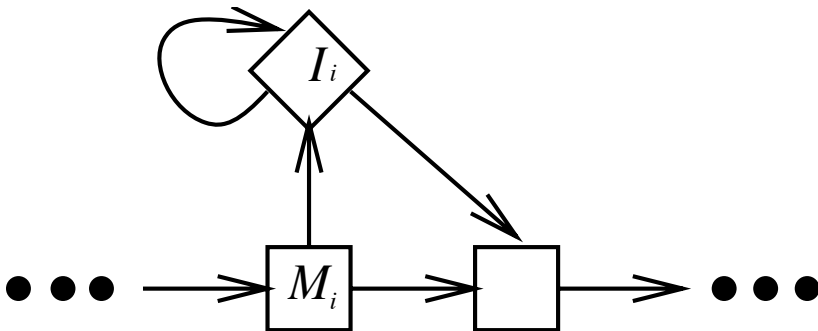
$$P(x \mid M) = \prod_{i=1}^L e_i(x_i)$$

Can, as usual, convert probabilities to log-odds score

Organization of a Profile HMM (2)

Insertion States

- But this assumes ungapped alignments!
- To handle gaps, consider insertions and deletions
 - Insertion: part of x that doesn't match anything in multiple alignment (use insert states)



Organization of a Profile HMM (3)

Deletion States

CSCE
471/871
Lecture 4:
Profile Hidden
Markov
Models

Stephen Scott

Organization

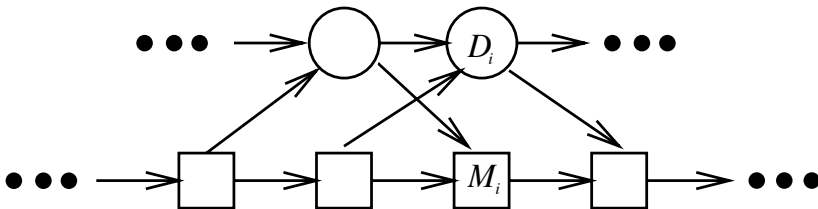
Match States
Insertion States

Deletion States
General Structure

Building
Models

Searching

- Deletion: parts of multiple alignment not matched by any residue in x (use silent delete states)



General Profile HMM Structure

CSCE

471/871

Lecture 4:

Profile Hidden

Markov

Models

Stephen Scott

Organization

Match States

Insertion States

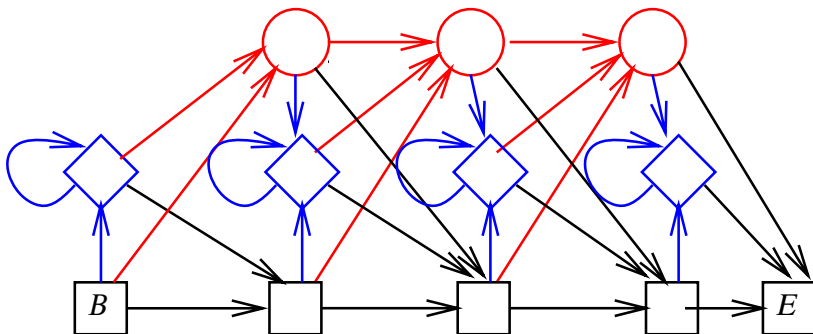
Deletion States

General Structure

Building

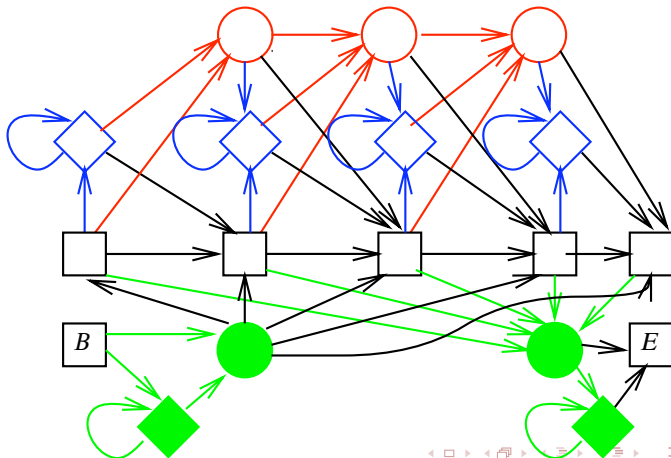
Models

Searching



Handling non-Global Alignments

- Original profile HMMs model entire sequence
- Add **flanking model states** (or **free insertion modules**) to generate non-local residues



Building a Model

Determining States

- Given a multiple alignment, how to build an HMM?
 - General structure defined, but how many match states?

CSCE
471/871
Lecture 4:
Profile Hidden
Markov
Models

Stephen Scott

Organization

Building
Models

States

Probabilities
Pseudocounts

Searching

...	V	G	A	-	-	H	A	G	E	Y	...
...	V	-	-	-	-	N	V	D	E	V	...
...	V	E	A	-	-	D	V	A	G	H	...
...	V	K	G	-	-	-	-	-	-	D	...
...	V	Y	S	-	-	T	Y	E	T	S	...
...	F	N	A	-	-	N	I	P	K	H	...
...	I	A	G	A	D	N	G	A	G	V	...

Building a Model (2)

Determining States

Given a multiple alignment, how to build an HMM?

- General structure defined, but how many match states?
- Heuristic: if more than half of characters in a column are non-gaps, include a match state for that column

...	V	G	A	-	-	H	A	G	E	Y	...
...	V	-	-	-	-	N	V	D	E	V	...
...	V	E	A	-	-	D	V	A	G	H	...
...	V	K	G	-	-	-	-	-	-	D	...
...	V	Y	S	-	-	T	Y	E	T	S	...
...	F	N	A	-	-	N	I	P	K	H	...
...	I	A	G	A	D	N	G	A	G	V	...

- Now, find parameters
- Multiple alignment + HMM structure \rightarrow state sequence

Diagram illustrating the alignment of a sequence with a reference sequence. The reference sequence is "V G A - - H A G E Y". The aligned sequence is "V - - - N V D E V". The alignment is shown by a red arrow pointing from the "A" in the reference sequence to the "A" in the aligned sequence. The alignment is also shown by a green arrow pointing from the "A" in the reference sequence to the "A" in the aligned sequence. The alignment is also shown by a green arrow pointing from the "A" in the reference sequence to the "A" in the aligned sequence.

Non-gap in match column -> match state

Gap in match column -> delete state

Non-gap in insert column -> insert state

Gap in insert column -> ignore

Durbin Fig 5.4, p. 109

Building a Model (4)

Estimating Probabilities

CSCE
471/871

Lecture 4:
Profile Hidden
Markov
Models

Stephen Scott

Organization

Building
Models
States

Probabilities

Pseudocounts

Searching

- Count number of transitions and emissions and compute:

$$a_{kl} = \frac{A_{kl}}{\sum_{l'} A_{kl'}}$$

$$e_k(b) = \frac{E_k(b)}{\sum_{b'} E_k(b')}$$

- Still need to beware of some counts = 0

Weighted Pseudocounts

CSCE
471/871

Lecture 4:
Profile Hidden
Markov
Models

Stephen Scott

Organization

Building
Models

States

Probabilities

Pseudocounts

Searching

- Let c_{ja} = observed count of residue a in position j of multiple alignment

$$e_{M_j}(a) = \frac{c_{ja} + Aq_a}{\sum_{a'} c_{ja'} + A}$$

- q_a = background probability of a , A = weight placed on pseudocounts (sometimes use $A \approx 20$)
- Background probabilities also called a prior distribution

CSCE

471/871

Lecture 4:

Profile Hidden

Markov

Models

Stephen Scott

Organization

Building

Models

States

Probabilities

Pseudocounts

Searching

- Can be thought of as a mixture of pseudocounts
- The mixture has different components, each representing a different context of a protein sequence
 - E.g., in parts of a sequence folded near protein's surface, more weight (higher q_a) can be given to hydrophilic residues
 - But in other regions, may want to give more weight to hydrophobic residues
- Will find a different mixture for each position of the alignment based on the distribution of residues in that column

Dirichlet Mixtures (2)

CSCE

471/871

Lecture 4:

Profile Hidden

Markov

Models

Stephen Scott

Organization

Building

Models

States

Probabilities

Pseudocounts

Searching

- Each component k consists of a vector of pseudocounts $\vec{\alpha}^k$ (so α_a^k corresponds to Aq_a) and a mixture coefficient (m_k , for now) that is the probability that component k is selected
- Pseudocount model k is the “correct” one with probability m_k
- We’ll set the mixture coefficients for each column based on which vectors best fit the residues in that column
 - E.g., first column of alignment on slide 10 is dominated by V, so any vector $\vec{\alpha}^k$ that favors V will get a higher m_k

Dirichlet Mixtures (3)

CSCE

471/871

Lecture 4:

Profile Hidden

Markov

Models

Stephen Scott

Organization

Building
Models

States

Probabilities

Pseudocounts

Searching

- Let \vec{c}_j be vector of counts in column j

$$e_{M_j}(a) = \sum_k P(k | \vec{c}_j) \frac{c_{ja} + \alpha_a^k}{\sum_{a'} (c_{ja'} + \alpha_{a'}^k)}$$

- $P(k | \vec{c}_j)$ are the posterior mixture coefficients, which are easily computed [Sjölander et al. 1996], yielding:

$$e_{M_j}(a) = \frac{X_a}{\sum_{a'} X_{a'}} ,$$

where

$$X_a = \sum_k m_{k0} \exp(\ln B(\vec{\alpha}^k + \vec{c}_j) - \ln B(\vec{\alpha}^k)) \frac{c_{ja} + \vec{\alpha}_a^k}{\sum_{a'} (c_{ja'} + \alpha_{a'}^k)}$$

$$\ln B(\vec{x}) = \sum_i \ln \Gamma(x_i) - \ln \Gamma\left(\sum_i x_i\right)$$

Dirichlet Mixtures (4)

- Γ is gamma function, and $\ln \Gamma$ is computed via `lgamma` and related functions in C
- m_{k0} is prior probability of component k ($= q$ below)

Parameters of Dirichlet mixture prior Blocks9

	Comp. 1	Comp. 2	Comp. 3	Comp. 4	Comp. 5	Comp. 6	Comp. 7	Comp. 8	Comp. 9
q	0.1829	0.0576	0.0898	0.0792	0.0831	0.0911	0.1159	0.0660	0.2340
$ \vec{\alpha} $	1.1806	1.3558	6.6643	2.0814	2.0810	2.5681	1.7660	4.9876	0.0995
A	0.2706	0.0214	0.5614	0.0701	0.0411	0.1156	0.0934	0.4521	0.0051
C	0.0398	0.0103	0.0454	0.0111	0.0147	0.0373	0.0047	0.1146	0.0040
D	0.0175	0.0117	0.4383	0.0194	0.0056	0.0124	0.3872	0.0624	0.0067
E	0.0164	0.0108	0.7641	0.0946	0.0102	0.0181	0.3478	0.1157	0.0061
F	0.0142	0.3856	0.0873	0.0131	0.1536	0.0517	0.0108	0.2842	0.0034
G	0.1319	0.0164	0.2591	0.0480	0.0077	0.0172	0.1058	0.1402	0.0169
H	0.0123	0.0761	0.2149	0.0770	0.0071	0.0049	0.0497	0.1003	0.0036
I	0.0225	0.0353	0.1459	0.0329	0.2996	0.7968	0.0149	0.5502	0.0021
K	0.0203	0.0139	0.7622	0.5766	0.0108	0.0170	0.0942	0.1439	0.0050
L	0.0307	0.0935	0.2473	0.0722	0.9994	0.2858	0.0277	0.7006	0.0059
M	0.0153	0.0220	0.1186	0.0282	0.2101	0.0758	0.0100	0.2765	0.0014
N	0.0482	0.0285	0.4415	0.0803	0.0061	0.0145	0.1878	0.1185	0.0041
P	0.0538	0.0130	0.1748	0.0376	0.0130	0.0150	0.0500	0.0974	0.0090
Q	0.0206	0.0230	0.5308	0.1850	0.0197	0.0113	0.1100	0.1266	0.0036
R	0.0236	0.0188	0.4655	0.5067	0.0145	0.0126	0.0386	0.1436	0.0065
S	0.2161	0.0291	0.5834	0.0737	0.0120	0.0275	0.1194	0.2789	0.0031
T	0.0654	0.0181	0.4455	0.0715	0.0357	0.0883	0.0658	0.3584	0.0036
V	0.0654	0.0361	0.2270	0.0425	0.1800	0.9443	0.0254	0.6617	0.0029
W	0.0037	0.0717	0.0295	0.0112	0.0127	0.0043	0.0032	0.0615	0.0027
Y	0.0096	0.4196	0.1210	0.0287	0.0264	0.0167	0.0187	0.1993	0.0026

Searching for Homologues

CSCE
471/871

Lecture 4:
Profile Hidden
Markov
Models

Stephen Scott

Organization

Building
Models

Searching

Viterbi
Forward
Aligning

Score a candidate match x by using log-odds:

- $P(x, \pi^* | M)$ is probability that x came from model M via most likely path π^*
 \Rightarrow Find using Viterbi
- $Pr(x | M)$ is probability that x came from model M summed over all possible paths
 \Rightarrow Find using forward algorithm
- $score(x) = \log(P(x | M)/P(x | \phi))$
 - ϕ is a “null model”, which is often the distribution of amino acids in the training set or AA distribution over each individual column
 - If x matches M much better than ϕ , then score is large and positive

Viterbi Equations

CSCE

471/871

Lecture 4:

Profile Hidden

Markov

Models

Stephen Scott

Organization

Building
Models

Searching

Viterbi

Forward

Aligning

- $V_j^M(i)$ = log-odds score of best path matching $x_{1\dots i}$ to model, x_i emitted by M_j (similarly define $V_j^I(i)$ and $V_j^D(i)$)
- B is M_0 , $V_0^M(0) = 0$, E is M_{L+1} (V_{L+1}^M = final)

$$V_j^M(i) = \log \left(\frac{e_{M_j}(x_i)}{q_{x_i}} \right) + \max \begin{cases} V_{j-1}^M(i-1) + \log a_{M_{j-1}M_j} \\ V_{j-1}^I(i-1) + \log a_{I_{j-1}M_j} \\ V_{j-1}^D(i-1) + \log a_{D_{j-1}M_j} \end{cases}$$

$$V_j^I(i) = \log \left(\frac{e_{I_j}(x_i)}{q_{x_i}} \right) + \max \begin{cases} V_j^M(i-1) + \log a_{M_jI_j} \\ V_j^I(i-1) + \log a_{I_jI_j} \\ V_j^D(i-1) + \log a_{D_jI_j} \end{cases}$$

$$V_j^D(i) = \max \begin{cases} V_{j-1}^M(i) + \log a_{M_{j-1}D_j} \\ V_{j-1}^I(i) + \log a_{I_{j-1}D_j} \\ V_{j-1}^D(i) + \log a_{D_{j-1}D_j} \end{cases}$$

Forward Equations

CSCE

471/871

Lecture 4:

Profile Hidden

Markov

Models

Stephen Scott

Organization

Building

Models

Searching

Viterbi

Forward

Aligning

$$F_j^M(i) = \log \left(\frac{e_{M_j}(x_i)}{q_{x_i}} \right) + \log [a_{M_{j-1}M_j} \exp(F_{j-1}^M(i-1)) + a_{I_{j-1}M_j} \exp(F_{j-1}^I(i-1)) + a_{D_{j-1}M_j} \exp(F_{j-1}^D(i-1))]$$

$$F_j^I(i) = \log \left(\frac{e_{I_j}(x_i)}{q_{x_i}} \right) + \log [a_{M_jI_j} \exp(F_j^M(i-1)) + a_{I_jI_j} \exp(F_j^I(i-1)) + a_{D_jI_j} \exp(F_j^D(i-1))]$$

$$F_j^D(i) = \log [a_{M_{j-1}D_j} \exp(F_{j-1}^M(i)) + a_{I_{j-1}D_j} \exp(F_{j-1}^I(i)) + a_{D_{j-1}D_j} \exp(F_{j-1}^D(i))]$$

$\exp(\cdot)$ needed for sums and logs (can still be fast; see p. 78)

Aligning a Sequence with a Model (Multiple Alignment)

CSCE

471/871

Lecture 4:

Profile Hidden

Markov

Models

Stephen Scott

Organization

Building

Models

Searching

Viterbi

Forward

Aligning

- Given a string x , use Viterbi to find most likely path π^* and use the state sequence as the alignment
- More detail in Durbin, Section 6.5
 - Also discusses building an initial multiple alignment and HMM simultaneously via Baum-Welch