

Profile HMMs for Sequence Families

COMP 571
Luay Nakhleh, Rice University

1

Sequence Families

- * Functional biological sequences typically come in families
- * Sequences in a family have diverged during evolution, but normally maintain the same or a related function
- * Thus, identifying that a sequence belongs to a family tells about its function

2

HMM Profile

- * Consensus modeling of the family using a probabilistic model
- * Built from a given multiple alignment (assumed to be correct)

3

4

Alignment of 7 globins

The 8 alpha helices are shown as A-H above the alignment

5

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

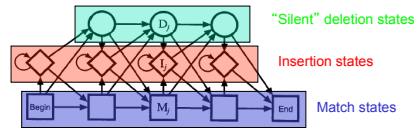
6

- $$S = \sum_{i=1}^L \log \frac{e_i(x_i)}{q_{x_i}}$$

ProfileHMMs - February 13, 2017

Adding Indels to Obtain a Profile HMM

7



Profile HMMs generalize pairwise alignment

Deriving Profile HMMs from Multiple Alignments

8

- * Essentially, we want to build a model representing the consensus sequence for a family, rather than the sequence of any particular member
- * Non-probabilistic profiles and profile HMMs

Non-probabilistic Profiles

9

- * Gribskov, McLachlan, and Eisenberg 1987
- * No underlying probabilistic model, but rather assigned position specific scores for each match state and gap penalty
- * The score for each consensus position is set to the average of the standard substitution scores from all the residues in the corresponding multiple sequence alignment column

Non-probabilistic Profiles

```

HBA_HUMAN  ...A--HAGEY...
HBB_HUMAN  ...V---NVDEV...
MYG_PHYCA  ...EA--DVAGH...
GLB3_CHITP  ...KG-----D...
GLB5_PETMA  ...YS--TYETS...
LGB2_LUPLU  ...HA--NTPKH...
GLB1_GLYDI  ...TAGADNGAGV...
          ...*
  
```

The score for residue 'a' in column 1

$$\frac{2}{3}s(V,a) + \frac{1}{3}s(F,a) + \frac{1}{3}s(I,a)$$

s(a,b) : standard substitution matrix

10

Non-probabilistic Profiles

- * They also set gap penalties for each column using a heuristic equation that decreases the cost of a gap according to the length of the longest gap observed in the multiple alignment spanning the column

11

Problem with the Approach

- * If we had an alignment with 100 sequences, all with a cysteine (C), at some position, the probability distribution for that column for an "average" profile would be exactly the same as would be derived from a single sequence
- * Doesn't correspond to our expectation that the likelihood of a cysteine should go up as we see more confirming examples

12

Similar Problem with Gaps

```
HBA_HUMAN   ...VGA--HAGEY...
HBB_HUMAN   ...V---HVDEV...
MYG_PHYCA   ...VEA--DVAGH...
GLB3_CHITP   ...VKG-----D...
GLB5_PETMA   ...VYS--TYETS...
LGB2_LUPLU   ...FNA--NTPKH...
GLB1_GLYDI   ...TAGADNGAGV...
            ***  *****
```

Scores for a deletion in columns 2 and 4 would be set to the same value

More reasonable to set the probability of a new gap opening to be higher in column 4

13

Basic Profile HMM Parameterization

- * A profile HMM defines a probability distribution over the whole space of sequences
- * The aim of parameterization is to make this distribution peak around members of the family
- * Parameters: probabilities and the length of the model

14

Model Length

- * A simple rule that works well in practice is that columns that are more than half gap characters should be modeled by inserts

15

Probability Values

$$a_{k\ell} = \frac{A_{k\ell}}{\sum_{\ell'} A_{k\ell'}} \quad e_k(a) = \frac{E_k(a)}{\sum_{a'} E_k(a')}$$

k, ℓ : indices over states

$a_{k\ell}, e_k(a)$: transition and emission probabilities

$A_{k\ell}, E_k(a)$: transition and emission frequencies

16

Problem with the Approach

- * Transitions and emissions that don't appear in the training data set would acquire zero probability (would never be allowed)
- * Solution: add pseudo-counts to the observed frequencies
- * Simple pseudo-count is Laplace's rule: add one to each frequency

17

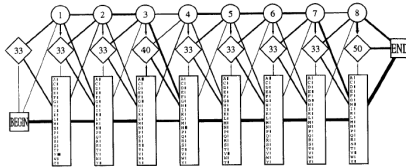
Example

```
HBA_HUMAN   ...VGA--HAGEY...
HBB_HUMAN   ...V---NVDEV...
MYG_PHYCA   ...VEA--DVAGH...
GLB3_CHITP   ...VRG-----D...
GLB5_PETMA   ...VYS--TVETS...
LGB2_LUPLU   ...FNA--NIPKH...
GLB1_GLYDI   ...IAGADNGAÖV...
***      *****
```

$e_{M_1}(V) = 6/27, e_{M_1}(I) = e_{M_1}(F) = 2/27$
 $e_{M_1}(a) = 1/27$ for all residue types a other than V, I, F
 $a_{M_1M_2} = 7/10, a_{M_1D_2} = 2/10$ and $a_{M_1I} = 1/10$

18

Example: Full Profile HMM



19

Searching with Profile HMMs

- * One of the main purposes of developing profile HMMs is to use them to detect potential membership in a family
- * We can either use Viterbi algorithm to get the most probable alignment or the forward algorithm to calculate the full probability of the sequence summed over all possible paths

20

Viterbi Algorithm

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

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

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

21

Forward Algorithm

$$\begin{aligned}F_j^M(i) &= \log \frac{e_{M_j}(x_i)}{q_n} + \log [a_{M_{j-1}M_j} \exp(F_{j-1}^M(i-1)) \\&\quad + a_{U_{j-1}M_j} \exp(F_{j-1}^I(i-1)) + a_{V_{j-1}M_j} \exp(F_{j-1}^D(i-1))]; \\F_j^I(i) &= \log \frac{e_{I_j}(x_i)}{q_n} + \log [a_{M_jI_j} \exp(F_j^M(i-1)) \\&\quad + a_{I_{j-1}I_j} \exp(F_{j-1}^I(i-1)) + a_{V_{j-1}I_j} \exp(F_{j-1}^D(i-1))]; \\F_j^D(i) &= \log [a_{M_{j-1}D_j} \exp(F_{j-1}^M(i)) + a_{U_{j-1}D_j} \exp(F_{j-1}^I(i)) \\&\quad + a_{V_{j-1}D_j} \exp(F_{j-1}^D(i))].\end{aligned}$$

22

Questions?

23
