CISC 636 Computational Biology & Bioinformatics (Fall 2016)

Hidden Markov Models (IV)

- a. Profile HMMs
- b. ipHMMs
- c. GeneScan
- d. TMMOD
- e. ipHMM and ETB-viterbi

Profile HMM for a family of sequences

Applications of HMM's

- Given a family of sequences, $O^{l}=O^{l}_{1}...O^{l}_{K'}$, build a hidden Markov model that best fits to this family-->Problem 3
 - Correct multiple alignment is given--> Problem 3, path known
 - MA built from structural information
 - MA obtained from other sequence based alignment procedures
 - Alignment is not assumed--> Problem 3, path not known (B-W)
- Use the obtained model to:
 - Score potential matches of new sequences-->Problem 1
 - Align new sequences--> Problem 2





Javier Garcia-Frias



CISC636, F16, Lec12, Liao



Profile HMM: Correct alignment assumed

Key idea of profile HMM

- Transition and emission probabilities capture specific information about each position in the multiple alignment of the whole family
- Profile HMM=Statistical model representing the family

Questions

- How do we build the profile HMM that best fits to a given family? -->Problem 3 (simplified)
- How do we detect potential membership in this family (for new sequences)? --> Problem 1
- How do we align a new sequence? --> Problem 2

Parameterization of profile HMM's: Correct alignment assumed

Profile HMM parametrization (simplified Problem 3)

- Model length
 - Length (and structure) completely defined when we decide which MA columns should be assigned to match states
 - Manual construction
 - Heuristic construction: e.g., column aligned if proportion of gaps is less than a threshold
 - More sophisticated methods
- Parameter estimation
 - Alignment is given-->Path through model is given for any sequence
 - Apply solution to Problem 3 when path is given (just count events)



Javier Garcia-Frias

Parameterization of profile HMM's: Correct alignment assumed

Emission probabilities: Estimate from number of emissions

$N(A M_1)=3$	$N(other M_1)=0$	I_0, I_1, I_3 are not used	
$N(A M_2)=3$	N(other M ₂)=0	$N(A I_2)=2$ $N(C I_2)=2$	$N(G I_2)=1$
$N(C M_3)=4$	N(other M ₃)=0	$N(L I_2)=1$ $N(V I_2)=1$	$N(other I_2)=0$

Transition probabilities: Estimate from number of transitions

 $N(M_1|B)=3 N(D_1|B)=1$ $N(M_2|M_1)=3 N(D_2|M_1)=1$ $N(M_3|M_2)=1 N(I_2|M_2)=2$

$$N(I_2|D_2)=1$$

$$N(I_2|I_2)=4$$
 $N(M_3|I_2)=3$

• If number of sequences is not high enough, estimation should be modified

Javier Garcia-Frias

 $N(E|M_3)=3$

Membership in a profile HMM

Detection of potential membership, for a new sequence, in family defined by a profile HMM (Problem 1)

- Apply forward equation
- Since $P(\mathbf{O}|M)$ is length dependent, usually scoring function is modified

Scoring=log
$$\frac{P(\mathbf{O}|M)}{P(\mathbf{O}|S)}$$

S is called "standard model": Model to use if sequences were independently distributed

• Other statistical approaches can also be used to improve the scoring system

Multiple alignment using profile HMM's

No alignment is assumed

- From an initially unaligned family of sequences, jointly perform:
 - Profile HMM estimation
 - Alignment estimation

1. Initialization

• Choose length of profile HMM and initialize parameters

2. Training

- Estimate parameters of the profile HMM
- Path not known (no alignment)--> Problem 3 (Baum-Welch)

3. Alignment

• Align all sequences using Viterbi algorithm (Problem 2)



Interaction profile HMM (ipHMM)



Friedrich et al, Bioinformatics 2006

GENSCAN (generalized HMMs)

- Chris Burge, PhD Thesis '97, Stanford
- <u>http://genes.mit.edu/GENSCAN.html</u>
- Four components
 - A vector π of initial probabilities
 - A matrix T of state transition probabilities
 - A set of length distribution f
 - A set of sequence generating models P
- Generalized HMMs:
 - at each state, emission is not symbols (or residues), rather, it is a fragment of sequence.
 - Modified viterbi algorithm CISC636, F16, Lec12, Liao



- Initial state probabilities
 - As frequency for each functional unit to occur in actual genomic data. E.g., as ~ 80% portion are non-coding intergenic regions, the initial probability for state N is 0.80
- Transition probabilities
- State length distributions

- Training data
 - 2.5 Mb human genomic sequences
 - 380 genes, 142 single-exon genes, 1492 exons and 1254 introns
 - 1619 cDNAs

Open areas for research

- Model building
 - Integration of domain knowledge, such as structural information, into profile HMMs
 - Meta learning?
- Biological mechanism DNA replication
- Hybrid models
 - Generalized HMM

— • • •

TMMOD: An improved hidden Markov model for predicting transmembrane topology





TMHMM by Krogh, A. et al JMB 305(2001)567-580



Accuracy of prediction for topology: 78%











Mod.	Reg.	Data set	Correct topology	Correct location	Sens- itivity	Speci- ficity
TMMOD 1	(a) (b) (c)	S-83	65 (78.3%) 51 (61.4%) 64 (77.1%)	67 (80.7%) 52 (62.7%) 65 (78.3%)	97.4% 71.3% 97.1%	97.4% 71.3% 97.1%
TMMOD 2	(a) (b) (c)	S-83	61 (73.5%) 54 (65.1%) 54 (65.1%)	65 (78.3%) 61 (73.5%) 66 (79.5%)	99.4% 93.8% 99.7%	97.4% 71.3% 97.1%
TMMOD 3	(a) (b) (c)	S-83	70 (84.3%) 64 (77.1%) 74 (89.2%)	71 (85.5%) 65 (78.3%) 74 (89.2%)	98.2% 95.3% 99.1%	97.4% 71.3% 97.1%
ТМНММ		S-83	64 (77.1%)	69 (83.1%)	96.2%	96.2%
PHDtm		S-83	(85.5%)	(88.0%)	98.8%	95.2%
TMMOD 1	(a) (b) (c)	S-160	117 (73.1%) 92 (57.5%) 117 (73.1%)	128 (80.0%) 103 (64.4%) 126 (78.8%)	97.4% 77.4% 96.1%	97.0% 80.8% 96.7%
TMMOD 2	(a) (b) (c)	S-160	120 (75.0%) 97 (60.6%) 118 (73.8%)	132 (82.5%) 121 (75.6%) 135 (84.4%)	98.4% 97.7% 98.4%	97.2% 95.6% 97.2%
TMMOD 3	(a) (b) (c)	S-160	120 (75.0%) 110 (68.8%) 135 (84.4%)	133 (83.1%) 124 (77.5%) 143 (89.4%)	97.8% 94.5% 98.3%	97.6% 98.1% 98.1%
TMHMM		S-160	123 (76.9%)	134 (83.8%)	97.1%	97.7%

Proteins Interact via Domains

Chemical bonds are formed between amino acids across interface at two interacting proteins. Do



Residues at interface tend to be more conserved due to selection pressure during evolution.

Domain A

Domain B

Profile Hidden Markov Models capturing interaction

- $P(x|\theta)$: probability that sequence x contains a domain described by the model θ .
- Viterbi algorithm can align x against the model to annotate interacting residues.



Capture long range correlation with HMM

(Kern, Gonzalez, Liao, Shanker, 2013)





$$S(x, y | x @h_i, y @h_{i+1}) = \log\left(\frac{p(y|x)}{p(x) p(y)}\right)$$









Fig. 5. A segment of the dynamic programming table for ipHMM. The lines show the path of an early traceback, with M_i states defined as hotspots. When hotspot (18, 17) is handled on the bottom right, traceback is performed until the previous hotspot, (6, 7), is reached.



Fig. 6. The improvement made by the ETB-Viterbi algorithm on the real data and the Top 1 simulated dataset.