

How to Score Multiple Alignments?

- Sum of Pairs Score (SP)
 - Optimal alignment: $O(d^N)$ [Dynamic Prog]
 - Approximate Algorithm: **Approx Ratio 2**
 - Locate Center: $O(d^2N^2)$
 - Locate Consensus: $O(d^2N^2)$

Consensus char: char with min distance sum

Consensus string: string of consensus char

Center: input string with min distance sum

Multiple Alignment Methods

- Phylogenetic Tree Alignment (**NP-Complete**)
 - Given tree, task is to label leaves with strings
- Iterative Method(s)
 - Build a MST using the distance function
- Clustering Methods
 - Hierarchical Clustering
 - K-Means Clustering

Multiple Alignment Methods (Cont'd)

- Gibbs Sampling Method
 - Lawrence, Altschul, Boguski, Liu, Neuwald, Winton, *Science*, 1993
- Hidden Markov Model
 - Krogh, Brown, Mian, Sjolander, Haussler, *JMB*, 1994

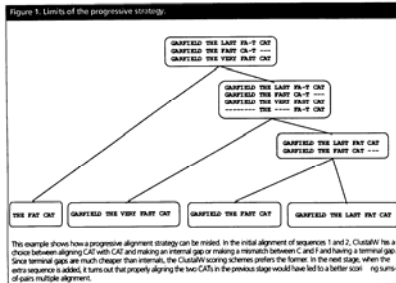
Multiple Sequence Alignments (MSA)

- Choice of Scoring Function
 - Global vs local
 - Gap penalties
 - Substitution matrices
 - Incorporating other information
 - Statistical Significance
- Computational Issues
 - Exact/heuristic/approximate algorithms for optimal MSA
 - Progressive/Iterative/DP
 - Iterative: Stochastic/Non-stochastic/Consistency-based
- Evaluating MSAs
 - Choice of good test sets or benchmarks (BAliBASE)
 - How to decide thresholds for good/bad alignments

Progressive MSA: CLUSTALW

REVIEW

Figure 1. Limits of the progressive strategy.



C. Notredame, *Pharmacogenomics*, 3(1), 2002.

Software for MSA

REVIEW

Table 1. Some recent and less recent available methods for MSA.			
MSA	Exact	http://www.bio.wustl.edu/bio/ma.html	198
MA	Iterative DCA	http://biochem.tch.fak.um-bielefeld.de/ma/	1911
Multalin	Progressive	http://www.toulouse.inra.fr/multalin.html	1913
ConsAlign	Consistency-based	http://www.gam.ac.uk/~ocap/m/	1798
Protein	Iterative/progressive	petrog@mm.mcg.ac.ca	1985
Protein	Iterative/Stochastic	http://ftp.genome.ucsf.edu/genomelab/ma-ccl/	1915
MSA-MER	Iterative/Stochastic/AM	http://primer.wustl.edu/	1988
GA	Iterative/Stochastic/GA	cchang@mathrow.uwaterloo.ca	1923

C. Notredame, *Pharmacogenomics*, 3(1), 2002.

MSA: Conclusions

- Very important
 - Phylogenetic analyses
 - Identify members of a family
 - Protein structure prediction
- No perfect methods
- Popular
 - Progressive methods: **CLUSTALW**
 - Recent interesting ones: **Prmp, SAGA, DiAlign, T-Coffee**
- Review of Methods [C. Notredame, *Pharmacogenomics*, 3(1), 2002]
 - **CLUSTALW** works reasonably well, in general
 - **DiAlign** is better for sequences with long insertions & deletions (indels)
 - **T-Coffee** is best available method

Profile Method

PROFILE METHOD, (M. Gribskov et al., '90)

Location in Seq.	Sequence	Protein Name
14	C V V A A A V	Ka RbqR
32	V V B H T I	Ec DwpR
33	V V P P T I	Ec RspD
76	A A L A T I	Ec YglR
138	C B R R V	Ec CAP
205	C L P R L	Ec AnaC
250	C L P R L	Sl AnaC
31	C V R E L I	Bt MerD

FREQUENCY TABLE

	A	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	Y
1	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	1	1	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	1	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
4	1	0	0	0	0	1	0	0	0	0	1	0	0	0	0	0	0	0	0	0
5	1	0	0	0	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
6	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	4	0	2	0	0	0	0	0	0	0	0	0	0	0	0

Profile Method

FREQUENCY TABLE

	A	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	Y
1	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	1	1	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
4	1	0	0	0	0	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0
5	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	4	0	2	0	0	0	0	0	0	0	0	0	0	0	0

WEIGHT MATRIX

	A	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	Y
1	104	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	21	79	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	21	0	0	0	0	79	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5	21	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6	21	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Weights(A,K) = log $\left(\frac{f_{AK}}{f_A f_K}\right) \cdot 100$

Profile Method

WEIGHT MATRIX

	A	C	R	G	T	X	L	M	F	R	D
1	0	104	0	104	0	0	0	0	0	0	0
2	21	78	0	0	0	0	44	0	0	0	0
3	0	0	23	0	0	0	46	0	0	0	103
4	21	0	32	0	38	32	0	0	0	86	39
5	21	0	42	23	0	0	0	0	0	0	72
6	21	0	0	0	0	0	0	0	0	0	63
7	0	0	0	0	0	98	0	44	0	0	0

Given the following protein sequence:

```

M T E D L F G D L Q D D T I L A H L D N
P A E D T S R F P A L L A E L N D L L R
G E L S R L G V D P A H S L E I V V A I
C K H L G G Q Q V T I P R G A L D S L
I R D L R I W N D F N G R N V S E L T T
R Y G V T F N T V Y K A I R R M R R L K
    
```

1/31/06

CAP5510/CGSS166

13

CpG Islands

- Regions in DNA sequences with increased occurrences of substring "CG"
- Rare: typically C gets methylated and then mutated into a T.
- Often around promoter or "start" regions of genes
- Few hundred to a few thousand bases long

1/31/06

CAP5510/CGSS166

14

Problem 1:

- **Input:** Small sequence S
- **Output:** Is S from a CpG island?
 - Build Markov models: M_+ and M_-
 - Then compare

1/31/06

CAP5510/CGSS166

15

Markov Models

+	A	C	G	T
A	0.180	0.274	0.426	0.120
C	0.171	0.368	0.274	0.188
G	0.161	0.339	0.375	0.125
T	0.079	0.355	0.384	0.182

-	A	C	G	T
A	0.300	0.205	0.285	0.210
C	0.322	0.298	0.078	0.302
G	0.248	0.246	0.298	0.208
T	0.177	0.239	0.292	0.292

1/31/06

CAP5510/CGSS166

16

How to distinguish?

- Compute

$$S(x) = \log\left(\frac{P(x|M+)}{P(x|M-)}\right) = \sum_{i=1}^L \log\left(\frac{p_{x(i-1)x}}{m_{x(i-1)x}}\right) = \sum_{i=1}^L F_{x(i-1)x}$$

r=p/m	A	C	G	T
A	-0.740	0.419	0.580	-0.803
C	-0.913	0.302	1.812	-0.685
G	-0.624	0.461	0.331	-0.730
T	-1.169	0.573	0.393	-0.679

Score(GCAC)
= .461-.913+.419
< 0.
GCAC not from CpG island.

Score(GCTC)
= .461-.685+.573
> 0.
GCTC from CpG island.

1/31/06

CAP5510/CGSS166

17

Problem 1:

- Input: Small sequence S
- Output: Is S from a CpG island?
 - Build Markov Models: M+ & M-
 - Then compare

Problem 2:

- Input: Long sequence S
- Output: Identify the CpG islands in S .
 - Markov models are inadequate.
 - Need Hidden Markov Models.

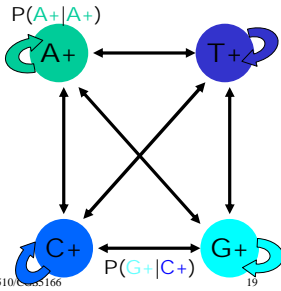
1/31/06

CAP5510/CGSS166

18

Markov Models

+	A	C	G	T
A	0.180	0.274	0.426	0.120
C	0.171	0.368	0.274	0.188
G	0.161	0.339	0.375	0.125
T	0.079	0.355	0.384	0.182



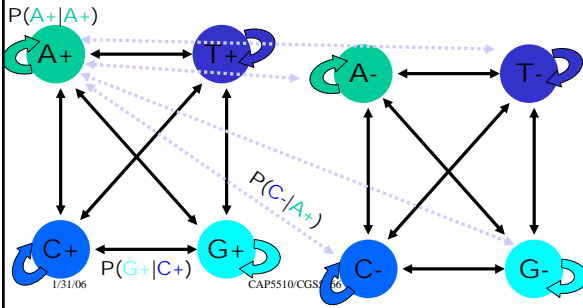
1/31/06

CAP5510/CGS5166

19

CpG Island + in an ocean of - First order Hidden Markov Model

MM=16, HMM= 64 transition probabilities (adjacent bp)



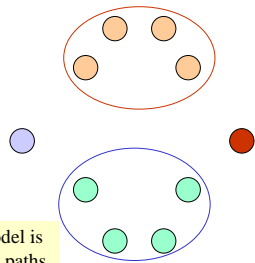
1/31/06

CAP5510/CGS5166

20

Hidden Markov Model (HMM)

- States
- Transitions
- Transition Probabilities
- Emissions
- Emission Probabilities



- What is **hidden** about HMMs?

Answer: The **path** through the model is hidden since there are many valid paths.

1/31/06

CAP5510/CGS5166

21

How to Solve Problem 2?

- Solve the following problem:

Input: Hidden Markov Model M ,
parameters Θ , emitted sequence S

Output: Most Probable Path Π

How: Viterbi's Algorithm (*Dynamic Programming*)

Define $\Pi[i,j]$ = MPP for first j characters of S ending in state i

Define $P[i,j]$ = Probability of $\Pi[i,j]$

- **Compute** state i with largest $P[i,j]$.

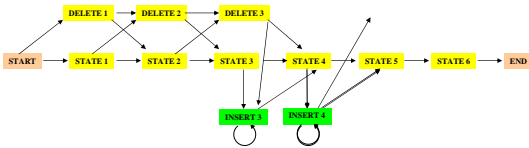
1/31/06

CAP5510/CGSS166

22

Profile HMMs with InDels

- Insertions
- Deletions
- Insertions & Deletions

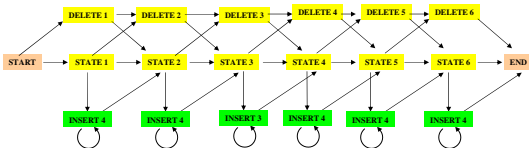


1/31/06

CAP5510/CGSS166

23

Profile HMMs with InDels



Missing transitions from DELETE j to INSERT j and
from INSERT j to DELETE $j+1$.

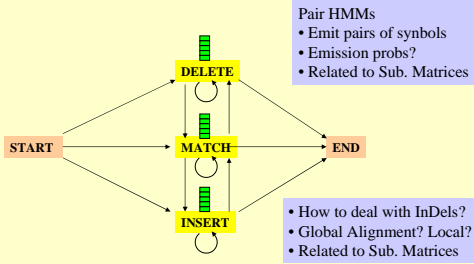
1/31/06

CAP5510/CGSS166

24

How to model Pairwise Sequence Alignment

LEAPVE
LAPVIE

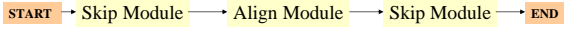


1/31/06

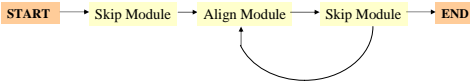
CAP5510/CGSS166

25

How to model Pairwise Local Alignments?



How to model Pairwise Local Alignments with gaps?



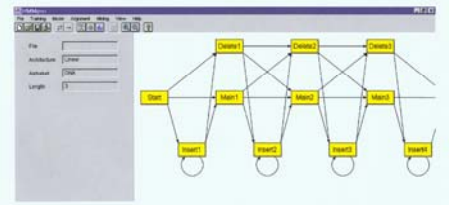
1/31/06

CAP5510/CGSS166

26

Standard HMM architectures

Linear Architecture



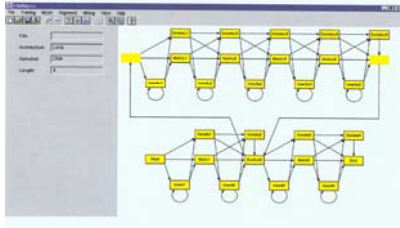
1/31/06

CAP5510/CGSS166

27

Standard HMM architectures

Loop Architecture



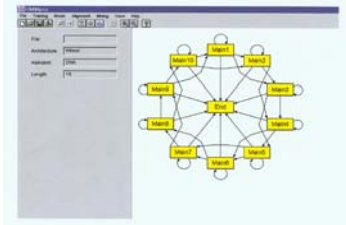
1/31/06

CAP5510/CGSS166

28

Standard HMM architectures

Wheel Architecture



1/31/06

CAP5510/CGSS166

29

Profile HMMs from Multiple Alignments

HBA_HUMAN VGA--HAGEY
HBB_HUMAN V----NVDEV
MYG_PHYCA VEA--DVAGH
GLB3_CHITP VKG-----D
GLB5_PETMA VYS--TYETS
LGB2_LUPLU FNA--NIPKH
GLB1_GLYDI IAGADNGAGV

Construct Profile HMM from above multiple alignment.

1/31/06

CAP5510/CGSS166

30

Iterative Solution to the **LEARNING QUESTION**
(Problem 5)

- Pick initial values for parameters Θ_0
- Repeat
 - Run training set S on model M
 - Count # of times transition $i \rightarrow j$ is made
 - Count # of times letter x is emitted from state i
 - Update parameters Θ
- Until (some stopping condition)

Entropy

- **Entropy** measures the variability observed in given data.

$$E = -\sum_c p_c \log p_c$$

- Entropy is useful in multiple alignments & profiles.
- Entropy is max when uncertainty is max.

G-Protein Couple Receptors

- Transmembrane proteins with 7 α -helices and 6 loops; many subfamilies
- Highly variable: 200-1200 aa in length, some have only 20% identity.
- [Baldi & Chauvin, '94] HMM for GPCRs
- HMM constructed with 430 match states (avg length of sequences); Training: with 142 sequences, 12 iterations

GPCR - Analysis

- Compute main state entropy values

$$H_i = -\sum_a e_{ia} \log e_{ia}$$

- For every sequence from test set (142) & random set (1600) & all SWISS-PROT proteins
 - Compute the negative log of probability of the most probable path π

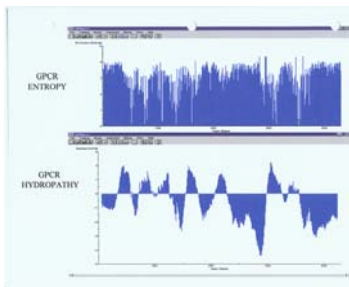
$$\text{Score}(S) = -\log(P(\pi | S, M))$$

1/31/06

CAP5510/CGSS166

37

GPCR Analysis



1/31/06

CAP5510/CGSS166

38

Entropy

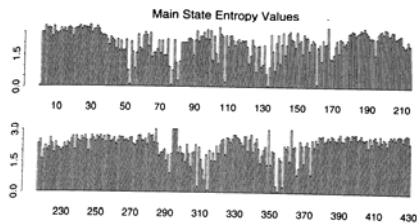


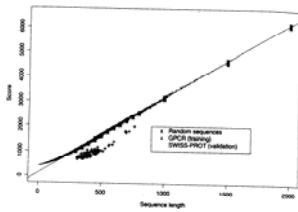
Figure 8.1: Entropy Profile of the Emission Probability Distributions Associated with the Main States of the HMM After 12 Cycles of Training.

1/31/06

CAP5510/CGSS166

39

GPCR Analysis (Cont'd)



1/31/06

CAP5510/CGSS166

40

Applications of HMM for GPCR

- Bacteriorhodopsin
 - Transmembrane protein with 7 domains
 - But it is not a GPCR
 - Compute score and discover that it is close to the regression line. Hence not a GPCR.
- Thyrotropin receptor precursors
 - All have long initial loop on INSERT STATE 20.
 - Also clustering possible based on distance to regression line.

1/31/06

CAP5510/CGSS166

41

HMMs – Advantages

- Sound statistical foundations
- Efficient learning algorithms
- Consistent treatment for insert/delete penalties for alignments in the form of locally learnable probabilities
- Capable of handling inputs of variable length
- Can be built in a modular & hierarchical fashion; can be combined into libraries.
- Wide variety of applications: Multiple Alignment, Data mining & classification, Structural Analysis, Pattern discovery, Gene prediction.

1/31/06

CAP5510/CGSS166

42

HMMs – Disadvantages

- Large # of parameters.
- Cannot express dependencies & correlations between hidden states.
