# Mathematical and Computational Molecular Biology
## Math/Biol 227 – Brendel – WQ94/95

**Time & Location:** Tu, Th 1:15P – 2:30P (Units: 3);  Gilbert 119
**Instructor:**  Volker Brendel (380-383A; 723-9256; volker@grendel.stanford.edu)

## Synopsis

Precipitated by an enormous increase in molecular sequence data (both DNA and protein), computational tools have become essential to molecular biology research. This course seeks to provide a state-of-the-art introduction to the subject. Emphasis will be on concepts and principles, combined with hands-on (-keyboard) applications.

## Prerequisites

This interdisciplinary course is directed at graduate and interested undergraduate students of biology, medicine, computer science, statistics, operations research, or mathematics. There are no formal prerequisites as all necessary knowledge will be developed in the course. However, some knowledge of the fundamental concepts of molecular biology and statistical analysis will be helpful. See the instructor for any questions regarding this.

# Mathematical and Computational Molecular Biology
## Math/Biol 227 – Brendel – WQ94/95

**Time & Location:**  Tu, Th 1:15P – 2:30P (Units: 3);     Gilbert 119
**Instructor:**             Volker Brendel (380-383A; 723-9256; volker@grendel.stanford.edu)

## Tentative Agenda

I **Tu, Jan. 10**   Overview.

II **Th, Jan. 12**   Pairwise sequence comparisons I; algorithms of Needleman & Wunsch and Smith & Waterman.

III **Tu, Jan. 17**   Pairwise sequence comparisons II; gap penalties; suboptimal alignments.

IV **Th, Jan. 19**   Score-based sequence analysis; single sequence features.

V **Tu, Jan. 24**   Score-based sequence analysis; pairwise sequence comparisons (SSPA).

VI **Th, Jan. 26**   Amino acid substitution scoring matrices.

VII **Tu, Jan. 31**   Query search methods; FASTA, BLAST.

VIII **Th, Feb. 2, S. Karlin instructor**   Sequence comparisons not requiring alignments.

IX **Tu, Feb. 7**   Phylogenetic trees from sequence data I.

X **Th, Feb. 9**   Phylogenetic trees from sequence data II.

XI **Tu, Feb. 14**   Runs, patterns, clusters of particular letter types.

XII **Th, Feb. 16, J. Kleffe instructor**   Word counts.

XIII **Tu, Feb. 21**   Spacings between sequence markers.

XIV **Th, Feb. 23**   Profile methods.

XV **Tu, Feb. 28**   Optimal signal search profiles (EM algorithm).

XVI **Th, March 2**   Hidden Markov chain Models.

XVII **Tu, March 7, D. Brutlag instructor**   Belief systems and neural networks for secondary structure prediction.

XVIII **Th, March 9**   Gene and intron prediction.

XIX **Tu, March 14**   Dead week lecture I.

XX **Th, March 16**   Dead week lecture II.

# Math & Computer Molecular Biology - Intro

I) Intro
 A) Sequences
  · DNA
  · protein

 B) Databases
  · Genbank, Genpept
  - EMBL, Swiss-prot

Nature 349:99 . Gilbert

II) Data-based
 Package-based
 Technique-based

III) Pairwise-Sequence Comparisons

a)
```
T T A C A G T T C  )              T T A C A G T T C
| | |  | | |  / /  } alignment    A T ▲ C A ▲ ▲ T C
A T C A T C
```

b) how score --- scoring matrix
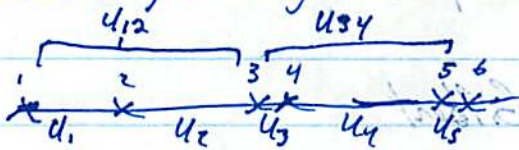
IV) Non-alignment based scores

V) Database
 - w/ lots of sequences in the database
   more likely getting a similarity by chance

VI) Profile
   - e.g REST. ENZYMES
   - expectation maximization ---finds motifs

VII) Distribution / Counting Words



(a) order $U_1 \cdots U_5$
(b) can also order $U_{12}$
                    $U_{34}$

_Math & Comp. Biology_

## SEQUENCE COMPARISON I - PAIRWISE

DNA $\downarrow$
$\{A, C, T, G\}_4$

PROTEIN $\downarrow$
$\{A, C, D, E, F, G, H, I, K, L, M, N, P, Q, R, S, T, V, W, Y\}_{20}$

### EXAMPLE

$A^7 = $ A G C C T A G
$B^7 = $ C A G C T G A
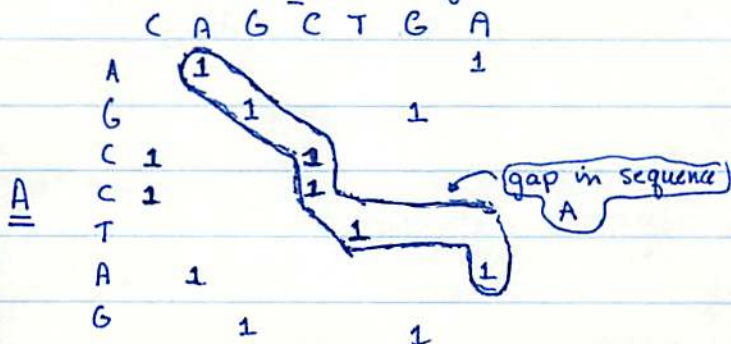
$\Big\}$

A G C C T A A G
C A G C A T G A

### ALIGNMENT PARAMETERS

① mismatches
② matches
③ gaps
④ overhangs

### PROBLEM - too many possibilities

### SOLUTION I - NEEDLEMAN & WUNSCH (1970) JMB 48: 443-453

Ⓐ represent all possible alignments on a matrix - as a path

can
- ~~must~~ have directionality
$\left(\begin{array}{c} e.g.\ 5'\ vs.\ 3' \\ N\ vs\ C \end{array}\right)$
∴ usually interested
in diagonal

$\underline{A}$

```
        B
     C  A  G  C  T  G  A
  A     1              1
  G     1  1     1
  C  1        1  1
  C  1        1  1     (gap in sequence)
  T              1      A
  A  1                 1
  G     1        1
```

- score for matches/mismatches
- alignments become diagonals
- gaps become vertical or horizontal lines

## NEEDLEMAN-WUNSCH II

- given two sequences

$$A_1^N = \{a_1, a_2 \cdots a_n\} \qquad B_1^M = \{b_1, b_2 \cdots b_m\}$$

Ⓐ an alignment is given by a set of paired indexes

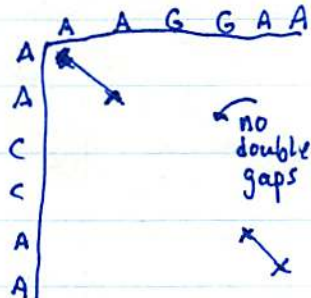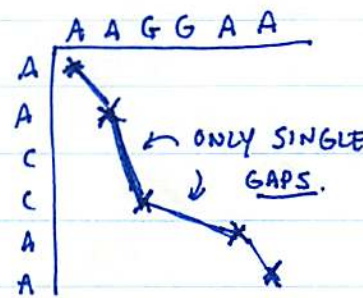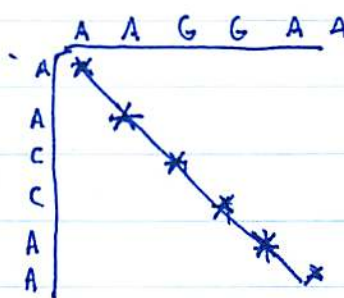- $(a_{j_1}, b_{k_1})(a_{j_2}, b_{k_2}) \cdots (a_{j_L}, b_{k_L})$

- with the restrictions

① $1 \leq j_1 < j_2 < \cdots j_L \leq N$ } always increasing
   $1 \leq k_1 < k_2 < \cdots k_L \leq N$

② if $j_L > j_{L-1} + 1$ then $k_L = k_{L-1} + 1$ } no double gaps
   $k_L > k_{L-1} + 1$ then $J_L = J_{L-1} + 1$

- that is in matrix the next match has to be in either next column or row.

· e.g   AACCAA   OR   AACG-AA   but not   AACC--AA
        AAGGAA        AA-GGAA             AA--GGAA



ONLY SINGLE GAPS.

no double gaps

Ⓑ SCORING

Score for
a path    $= S_\rho = S_\rho(A, B)$

substitution scores

weight for gaps $w(L) \leq 0 = \text{cost}$

# of gaps in each gap in A

gaps in each gap in B

$$= \sum_{\ell} \sigma(a_{j_\ell}, b_{k_\ell}) + \sum_{L} \omega(J_L - J_{L-1} - 1) + \sum_{L} \omega(k_L - k_{L-1} - 1)$$

- substitution scores
   Ⓐ SIMPLE

$$\sigma(a, b_i) = 1 \quad \text{if} \quad a = b \qquad \omega(k) = -K$$
$$\sigma(a, b) = 0 \quad \text{if} \quad a \neq b$$

© OPTIMIZATION

① FIND MAX $S_p$ FOR ALL PATHS

② ALGORITHM FOR FINDING MAX.  — depends on "no double gap" requirement

Ⓐ SET $S_{\emptyset j} = \omega_j$ for all $j = 1, 2 \cdots M$ } sets gap penalties
$S_{i\emptyset} = \omega_i$ for all $i = 1, 2 \cdots M$

optimal score for aligning $A_1^i = a_1, \cdots a_i$
$B_1^j = B_1 \cdots B_j$

$$S_{ij} = \max \begin{cases} S_{i-1, j-1} + \sigma(a_i, b_j) & A \\ S_{i, j-k} + \omega_k \text{ for all } k = 1, 2 \cdots j & B \\ S_{i-k, j} + \omega_k \text{ for all } k = 1, 2 \cdots i & C \end{cases}$$

then $S = S_{NM}$

PROOF

- optimal score for aligning $\mathcal{B} i, j$ is

~~F/ MA(i, j)algined~~

Ⓐ optimal score for aligning everything before $i, j$ plus score of aligning $i, j$

$= S_{i-1, j-1} + \sigma(a_i, b_j)$

$$\overline{\phantom{==}} \!\! {}^i_j$$

these are the only three possibilities to align $i$ & $j$

Ⓑ $\overline{\phantom{===}} {}^{i \cdots}_{\cdots j}$

$= S_{i, j-k} + \omega_k$

If assume no double gaps

Ⓒ $\overline{\phantom{===}} {}_{j \cdots}^{\cdots j}$

$= S_{i-k, j} + \omega_k$

RUNNING NW

```
         Δ  C  A  G  C  T  G  A
    Δ  0  -1  -2  -3  -4  -5  -6  -7
    A -1   0 → 0   □
    G -2
    C -3
    C -4
    T -5
    A -6
    G -7
```
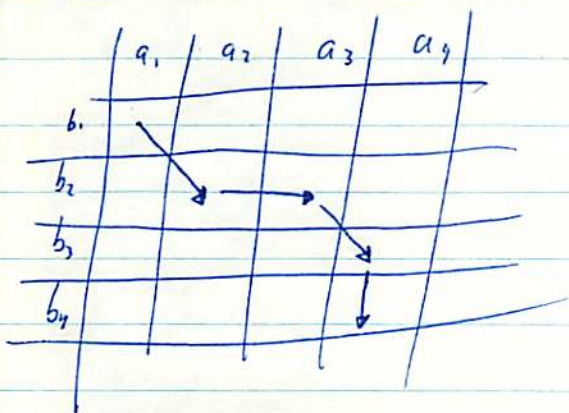
Score for each cell = <u>Max</u>
of all 3 adjacent cells

```
      A
  C   A
```

Mathematic & Computation Molecular Biology Pairwise Alignments

$A = \{a_1, a_2 \ldots a_m\}$

$B = \{b_1, b_2 \ldots b_n\}$

alignment = association of each $a_1 \ldots a_m$ w/ each $b_1 \ldots b_n$ (+ gaps)



$\nearrow$ = assoc. of $a_A$ $b_c$

$\downarrow$ = gap in b

$\rightarrow$ = gap in a

Needleman-Wunsch - corrected
 - no right angles  *
 - no double gaps

$P = path = \{a_{i_1}, b_{j_1} \ldots a_{i_\ell}, b_{j_\ell}\}$

$0 = i_0 < i_1 < \ldots < i_\ell < i_{\ell+1} = m+1$      $\ell = \#$ of diagonals = alignment of positions

$S_P = \underset{\text{for path}}{\underline{Score}} = \sum_{k=1}^{\ell} \sigma \{a_{i_k}, b_{j_k}\} + \sum_{i=1} gap(k)$

$S = _{max\ score} = Max\ S_P(A, B)$

b

a



$i,j$

$S_{ij} = \max \{$

— for last position, multiple paths

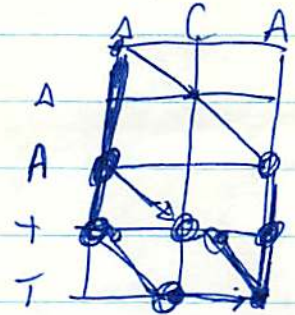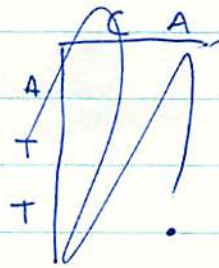if $\underline{a}$    $S = S_{i-1, j-1} + \sigma \{a_i, b_j\}$

if $b_k$    $S = S_{i-1, j-k-1} + \sigma(a_i, b_{j-k}) + \omega_k$

if $C_k$    $S = S_{i-k, j-1} + \sigma(a_{i-k}, b_j) + \omega_k$

        $S = S_{i-1-k, j-1}$

$-2 + 10 = 8$

Example    match = 10

           mismatch = * -10

           $\omega_1 = -2$

           $\omega_2 = -5$

           $\omega_3 = -8$



$-5 - 10 = -15 -2$

$= -17$

$-2$

| OLD NW | | Δ | C | A | | A | C | A |
|---|---|---|---|---|---|---|---|---|
| | Δ | 0 | -2 | -5 | | Δ | | |
| | A | -2 | -4 | ⊕ | | A | | |
| | T | -5 | -6 | 6 | | T | | |
| | T | -8 | 8 | 4 | | T | | |

— can only work
if $k\omega_1 \leq \omega_k$
— that is gap cannot
jump too much

— this # is
wrong — can
never get
4

<u>GOTOH</u> Let $w_k = -\alpha - \beta K$    $K = 1, 2...$    $\alpha, \beta > 0$    <span style="float:right">AFFINE GAP<br>PENALTIES</span>

$\therefore$ for all $k, \alpha, \beta$ $k w_1 \lesssim w_k$

· This allows speeding up of algorithm



$$H_{IJ} = \max \left\{ \begin{array}{l} S_{I, J-1} + w_1 \\ S_{I, J-k} + w_k \quad k=2...n \end{array} \right.$$

$$= \max \left\{ S_{I, J-1} \right\}$$

Proof

$$H_{IJ} = \max \left\{ S_{I, J-K} + w_k \right\} \; k = 1...n \; = \; \max \left\{ \begin{array}{l} S_{I, J-1} + w_1 \\ H_{I, J-1} - \beta \end{array} \right\}$$

$$V_{IJ} = \max \left\{ S_{J-k, J} + w_k \right\} \; k = 1...n$$

$$S_{ij} = \max \left\{ \begin{array}{l} S_{I-1, J-1} + \sigma(a_I, b_J) \\ H_{IJ} \\ V_{IJ} \end{array} \right.$$

· These assume that everthing ends in corner

· Thus end gap penalties are assessed

```
A  A T T A C
-  A T T - -
```

· <u>But</u> what about if you didn't want to assess these penalties

I) <u>MODIFICATION</u> <u>I</u>

  - start w/ all O's in columns    $S_{I0} = S_{0J} = 0$

  - find <u>max</u> score in <u>last</u> row or column

  $$S = \max \left\{ S_{MJ}, S_{IN} \right\}$$

# MODIFICATION II - LOCAL ALIGNMENT

A A [A T T A C] A A A
[A T T G C] G G G

$S_{i0} = S_{0j} = 0$   a) MAKE ~~1st~~ △ COLUMN/ROW ALL 0

b) AS SOON AS YOU GET A NEGATIVE # ... MAKE IT $\emptyset$.

$$S_{ij} \cancel{max} = \max \begin{cases} S_{i-1,j-1} + \sigma(a_i, b_j) \\ H_{i,j} \\ V_{i,j} \\ \emptyset \end{cases}$$

# Mathematical & Computational Molecular Biology

## FUNDAMENTALS

① alignment ⟷ lattice path

② scoring

- substitutions

DNA = 10 scores    $4 \times 4$ matrix $= 4 + 3 + 2 + 1$

protein = 210 scores    $20 \times 20$ matrix $= 20 + 19 + \dots + 2 + 1$

- gap penalties

- most use affine gap penalties    $\alpha + \beta k = w_k$ ; $\alpha, \beta < 0$

③ algorithm - how find high scoring alignments

④ interpretation

- how meaningful are alignments (statistics)

## LOCAL - SMITH - WATTERMAN

- CATTGC vs ATG
- mismatch = -1
- match = +3

$w_k = -1 - k$



|   | C | A | T | T | G | C |
|---|---|---|---|---|---|---|
|   | 0 ⋯ 0 | 0 | 0 | 0 | 0 | 0 |
| A | 0 | 0 | ③ | ① | 0 | 0 | 0 |
| T | 0 | 0 | 1 | ⑥ | ④ ⋯ 3 ⋯ 2 |
| G | 0 | 0 | | 4 | 5 | ⑦ | 5 |

multiple pathways

① take max $\left\{ \begin{array}{l} \text{all poss. scores} \\ 0 \end{array} \right\}$

② 0's in first row

③ fill #'s

④ pick highest score

⑤ trace path
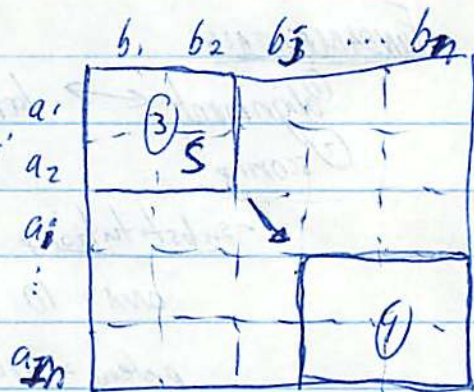
<u>So</u> - what do you do with <u>multiple paths</u>?

<u>Zuker</u> - Suboptimal Alignments

① constrain alignments such that they must go through certain points $a_I', b_J'$

② find all $(a_I, b_J)$ such that there is at least one alignment containing $(a_I', b_J')$ with score $S \geq S_{MN} - \delta$



two separate subproblems for each point
③ oi want to find $S_{I-1, J-1}$ = optimal score

④ then solve for this region by <u>inverting</u> sequences $(S_{I+1 \cdots M, J+1 \cdots N})$

brutlagxengen

<u>Nbor & Brutlag</u> - Near optimal alignment
- motivation ... are structural & optimal alignments the same
- Vimbroh & Argos ... reliably aligned regions

<u>Very parameters</u>
Dayhoff = PAM matrix

why not plot $w/$ <u>score</u> in 3D on Z axis

## Math/Comp

username @ grendel  /USER2/username

## SEQF

<u>ZUKER</u> - allows double gaps

$$\frac{aaa}{} \frac{bbb}{} \quad w_{k+l} \text{ <u>not</u> } w_k + w_l$$

## STATISTICS

(a) for full NW alignment w/ gaps --- no good theory for getting statistics
  - so (a) shuffle the seqs 100x · 1000x ··
    (b) order the scores
    (c) how many is score <u>above</u>

- Problem
  (a) computer intensive
  (b) changes w/ gap penalties ··· must redo for each

- Solutions -
  (1) <u>Disallow gaps</u> (<u>or</u> v.v.high penalty)
    (a) <u>all</u> alignments on diagonals
    (b) take a particular diagonal
      $$b_1 \cdots b_n$$
      $$a_1 \cdots a_{1+n-i}$$
    (c) replace alignment w/ sequence it scores
    (d) find segment <u>pairs</u> w/ highest aggregate score

## Statistical Theory

- Scores : $S_1, S_2, \ldots S_r$   $r$ different scores

w/ probabilities: $p_1, p_2, \ldots p_r$

random variable   $x_1, x_2, \ldots x_N$   prob $\{x_i = S_k\} = p_k$

$$S_k = \sum_{i=1}^{k} x_i$$

$$M = \max \{S_\ell - S_k\} \quad 0 \leq k \leq \ell \leq N$$

### EXAMPLE

| | $p_k$ |
|---|---|
| $S_1 = -5$ | 0.2 |
| $S_2 = -3$ | 0.2 |
| $S_3 = -1$ | 0.2 |
| $S_4 = 1$ | 0.2 |
| $S_5 = 5$ | 0.2 |

$$\overbrace{-3, \; -1, \; 5, \; 1, \; 5, \; 5, \; -3}^{16}$$

$S_1 = -3$    $S_4 = 2$    $S_7 = 9$

$S_2 = -4$    $S_5 = 7$

$S_3 = 1$     $S_6 = 12$

$$Max = S_6 - S_2 = 16$$

### 2 Assumptions

$$E(X) = \sum_{k=1}^{r} p_k S_k < 0$$

$$Prob(x_1 > 0) > 0$$

### RESULTS



want maximum increase

$E_k > 0$

$S_k$ vs $k$

EXCURSION PLOT



$E_k$

$E_k = S_p$

$k$

$\underbrace{\quad}_{M}$

$S_p = $ threshold

$E_0 = 0$

$E_k = \max \begin{cases} E_{k-1} + S_{k} \\ \text{OR} \\ \emptyset \end{cases}$

So... how calculate $S_p$?

$$\text{Prob}\left\{M > \frac{\ln N}{\lambda} + x\right\} \cong 1 - e^{-ke^{-\lambda x}}$$

$N$ = length   $\lambda$ } positive
$x$ = pos. variable   $k$ } parameters

$$\text{Prob}\left\{M > S_p\right\} \cong 1 - e^{-kN e^{-\lambda S_p}}$$

① $\lambda$ is unique positive root for $\sum_{k=1}^{n} P_R e^{\lambda S_k} = 1$

$\lambda$ is scale parameter

Of have scores $S_1, S_2 \cdots S_r$
w/ $\lambda, S_p$

$\Downarrow$

$e^{\lambda S_k} = e^{\lambda' \alpha S_k}$

Of have other-scores $S_{s2} \cdots S_{br} = \alpha S_k$

② $\lambda' \alpha = \lambda$   $S_p' = \alpha S_p$
$\lambda' = \lambda/\alpha$



$\sum_{k=1}^{n} P_R e^{\lambda S K}$

want this

② $k = \dfrac{F e^{-\lambda(A+B)}}{\lambda C}$

$$A = \sum_{k=1}^{\infty} \frac{1}{k} E\left[e^{\lambda S_k}; S_k < 0\right]$$

$$B = \sum_{k=1}^{\infty} \frac{1}{k} \text{prob}\left\{S_k \geq 0\right\}$$

$$C = E\left[S_1 e^{\lambda S_1}\right]$$

$F$ = correction factor for non-additive

# Mathematical Molecular Biology

mm bjc@cmgn ... gbrand

SEQUENCE OF SCORES $\quad X_1, X_2 \cdots X_n \qquad c$

prob. $\{X_i = S_k\} = P_k \qquad c = 1 \cdots N \quad k = 1 \cdots r$

$$S_k = \begin{cases} 0 \end{cases}$$

$P_k$ = frequency over all the sequence

$$M_{ax} = \max \{ S_e - S_k \}$$

$$PROB\{M = S\} = 1 - e^{-kNe^{-\lambda S}}$$

## APPLICATION

- for given scheme & significance level $p$
want $S_p$ where prob $\{M > S_p\} = p$

$$S_p = \frac{1}{\lambda}(\ln N + \ln K - \ln(\ln(1-p)))$$

SAPS
- anonymous FTP
on a gnomic
- Karlin-c
  calculates $\lambda, k$

## COMPOSITION BIAS

- in high scoring segments the occurrence frequences of scores are biased

$$q_k = P_k e^{\lambda S_k} \qquad \text{①} \sum q_k = 1 \qquad \text{②} \text{ if } S_k > 0 \Rightarrow q_k > P_k$$
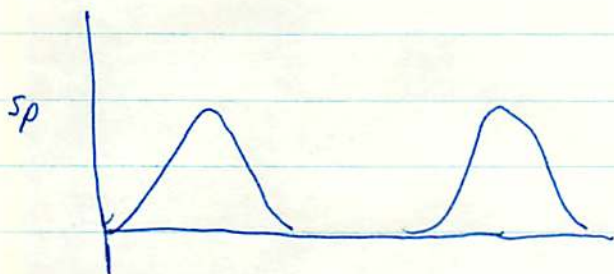$$\text{if } S_k < 0 \Rightarrow q_k < P_k$$

$S_k \propto \ln \frac{q_k}{P_k}$

- this can be reversed to calculate $S_k$ for different biases in alignments of interest

## EXPECTED LENGTH $\quad E[L] = \frac{S_p}{\sum q_k S_k}$

# APPLICATIONS - TRANSMEMBRANE DOMAINS

① identify region
② make g_k
③ evaluate w/ excursion plot



_but_ might miss segments containing two distinct domains

# APPLICATION - CHARGE CLUSTERS

|        | slay |      |
|--------|------|------|
| K, R   | +2   | +2   |
| D, E   | -2   | -8   |
| others | -1   | -5   |
|        | cluster | run of + chgs |

# HYDROPHOBICITY

# 2ARY STRUCTURE

doesn't work well because no positional information

# Application to Sequence Comparisons
   — for a given alignment

$M = $ length $\qquad b_1 \quad b_2 \; \text{---} \; b_{N-i+2} \qquad\qquad \leftarrow p'_k \qquad k = 1 \cdots n$

$N = $ length $\qquad a_i \quad a_{i+1} \; \text{.....} \; a_N \qquad\qquad\quad \leftarrow p_k \qquad k = 1 \cdots n$

$\qquad\qquad\qquad \sigma \quad \sigma'$

· sequence of scores $\quad \sigma \qquad \sigma \qquad \sigma$

· So get scores for all possible diagonals (no gaps)

· $\underline{M} = $ maximal segment pair (aggregate score)

$\text{Prob} \{ M > S \} = 1 - e^{-kMNe^{-\lambda S}}$ 
   — $k, \lambda$ analogous to before
   — $\lambda = $ unique positive root

Target frequencies = freq. of substitution

$$q_{ij} = p_i \, p_j \, e^{\lambda S_{ij}}$$

$$S_{ij} \propto \frac{q_{ij}}{p_i \, p_j}$$

$p_k$ must be somewhat similar to $p'_k$

## Mathematical & Computational Molecular Biology

### Amino Acid Substitution Scores

① 20 × 20 matrix
- almost always symmetrical

② given these scores ... high scoring segment pairs have a ~~~~ biased composition

$$q_{ij} = P_i P_j e^{\lambda S_{ij}}$$

prob. of generating
sequences i, j

↳ scale factor

$$S_{IJ} = \rho \ln \frac{q_{IJ}}{P_I P_J}$$

↳ proportionality constant

→ log-odds score

= these scores will target those regions w/ $q_{IJ}$

③ Assumptions
ⓐ at least one score > 0
ⓑ E $\sum P_i S_i < 0$

### Derivation of Substitution Scores

ⓐ A C L L M A G
   A C V I M G A

ⓑ count substitutions $k_{ij}, k_{ji}$

ⓒ symmetrize $C_{ij} = k_{ij} + k_{ji}$

ⓓ take row, column sums

rows = $C_{1\cdot}, C_{2\cdot}, C_{3\cdot} \cdots C_{20\cdot}$

columns = $C_{\cdot 1}, C_{\cdot 2}, C_{\cdot 3} \cdots C_{\cdot 20}$

$$C_{IJ} = \# \text{ of subs. } I \leftrightarrow J$$
$$C_{II} = 2\times \# \text{ of } I \leftrightarrow I \text{ matches}$$
$$C_{\cdot I} = C_{I\cdot} = \# \text{ of } I \text{ residues}$$
$$C_{\cdot\cdot} = \text{total } \# \text{ residues}$$

| $C_{11}$ | $C_{12}$ | $\cdots$ | $C_{1,20}$ |
|---|---|---|---|
| $C_{21}$ | $C_{22}$ | | |
| $C_{i\cdot}$ | | | |
| $\vdots$ | | | |
| $C_{20,}$ | $C_{\cdot,c}$ | | |

$$f_i = \frac{c_{i\cdot}}{c_{\cdot\cdot}} = \text{freq of residue } i$$

(3) Getting log-odds scores

    (a) Contingency Table Approach

$$S_{IJ} = \ln \frac{C_{IJ}}{(C_{I\cdot})(C_{\cdot J})/C_{\cdot\cdot}} = \frac{obs}{expected}$$

$$= \ln \frac{f_{IJ}}{f_I \cdot f_J}$$

(4) How get counts?

    (a) Blosum (Henikoff & Henikoff)

        - take blocks

        - each position make a column vector of diff residues

        - consider all pairwise comparisons

        - each used for substitutions

| | | | | | A | C | S | |
|---|---|---|---|---|---|---|---|---|
| 5A | | AA = 10 | | AC = 15 | A | 20 | 15 | 10 | 45 |
| 3C | | CC = 3 | | CS = 6 | C | 15 | 6 | 6 | 27 |
| 2S | | SS = 1 | | AS = 10 | S | 10 | 6 | 2 | 18 |
| | | | | | | 45 | 27 | 18 |

(b) STRUCTURE

(c) Pw Alignments

(d) PAMs - inferred from evolutionary data

① PAMs Dayhoff

  ⓐ infer substitutions from a _likely_ evolutionary _tree_
          - only works for highly conserved proteins

  ⓑ example - 4 seqs

  ACGH     DBGH     ADIJ    CBIJ
  2 B→C  \    / A→D   B→0  \    /
      ABGH              ABIJ
           \            /
              ABGJ

  ⓒ use this to get substitution counts $k_{IJ}$
  ⓓ symmetrize to $c_{IJ}$ — so both directions the same
→ ⓔ use transition probabilities of _change_ over _time_
      - _MARKOV MODEL_

ASSUMPTIONS
- mutations occur
    at constant rate
- accepted changes

      - changes are governed by probability transition matrices

      $\underset{(20 \times 20)}{M_{IJ}}$ = prob. that J changes to I in 1 unit of time

          $\sum_{i=1}^{20} M_{IJ} = 1$

PAMS

① Define relative mutability — how often does a particular aa change?

$$M_J = \frac{C_{\cdot J} - C_{JJ}}{C_{\cdot J}} = \frac{\text{\# residue involved in substitution}}{\text{all of that residue}}$$

② 1-step transition probability → prob. that J changes

(a) $\underline{M_{JJ} = 1 - \rho\, m_J}$

— $\rho \neq 1\bar{3}$ set in way such that
\# changes in 1 unit time = 1%

set $\sum\limits_{J=1}^{20} f_J\, M_{JJ} = 0.99$ = fraction not changing

Solve $\rho = \dfrac{1}{100 \sum\limits_{J=1}^{20} f_J m_J}$     1 accepted point mutation per 100 residues

③ $I \neq J$

$$\underline{M_{IJ} = \rho\, m_J \dfrac{C_{IJ}}{C_{\cdot J} - C_{JJ}}}$$

PROPERTIES

① $\sum\limits_{I=1}^{20} M_{IJ} = 1$

② $M_{IJ}^{\,n} = $ prob. of $J \to I$ in $n$ units

$$M_{IJ}^{(n)} = (M_{IJ})^n$$

③ Stationary distribution

$$\therefore (M_{IJ}) \begin{pmatrix} f_1 \\ f_2 \\ \vdots \\ f_{20} \end{pmatrix} = \begin{pmatrix} f_1 \\ f_2 \\ \vdots \\ f_{20} \end{pmatrix}$$

— ∴ only works for proteins with the same composition

④ $f_J M_{IJ} = f_I M_{JI}$
$f_J M_{IJ}^{\,n} = f_I M_{JI}^{\,n}$     } ∴ SYMMETRICAL LOG ODDS SCORES

$$\text{(8)} \quad S_{IJ}^{\,n} = \ln \frac{M_{IJ}^{\,n}}{f_I}$$

LOG-ODDS SCORES

as $n \to \infty$ $M_{IJ} \to f_I$

$\therefore$ all scores $\to \emptyset$

## Which PAM's

$^0\underline{n}$ small $\cdots$ for highly conserved

$n$ large $\cdots$ for highly divergent

## References: Score-based sequence analysis.

Altschul, S.F. (1991). Amino acid substitution matrices from an information theoretic perspective. *J. Mol. Biol.* **219**, 555–565.

Altschul, S.F. (1993). A protein alignment scoring system sensitive at all evolutionary distances. *J. Mol. Evol.* **36**, 290–300.

Collins, J.F., Coulson, A.F.W. & Lyall, A. (1988). The significance of protein sequence similarities. *Comput. Appl. Biosci.* **4**, 67–71.

Dayhoff, M.O., Schwartz, R.M. & Orcutt, B.C. (1978). A model of evolutionary change in proteins. In *Atlas of Protein Sequence and Structure* (Dayhoff, M.O., ed.). vol. 5, suppl. 3, pp. 345–352. Nat. Biomed. Res. Found., Washington, DC.

Dembo, A. & Karlin, S. (1993). Central limit theorems of partial sums for large segmental values. *Stoch. Proc. Appl.* **45**, 259–271.

Henikoff, S. & Henikoff, J.G. (1992). Amino acid substitution matrices from protein blocks. *Proc. Natl. Acad. Sci. U.S.A.* **89**, 10915–10919.

Henikoff, S. & Henikoff, J.G. (1993). Performance evaluation of amino acid substitution matrices. *Proteins* **17**, 49–61.

Jones, D.T., Taylor, W.R. & Thornton, J.M. (1992). The rapid generation of mutation data matrices from protein sequences. *Comput. Appl. Bosci.* **8**, 275–282.

Karlin, S. (1994). Statistical studies of biomolecular sequences: score-based methods. *Phil. Trans. R. Soc. Lond.* B **344**, 391–402.

Karlin, S. & Altschul, S.F. (1990). Methods for assessing the statistical significance of molecular sequence features by using general scoring schemes. *Proc. Natl. Acad. Sci. U.S.A.* **87**, 2264–2268.

Karlin, S. & Brendel, V. (1992). Chance and statistical significance in protein and DNA sequence analysis. *Science* **257**, 39–49.

Karlin, S. & Cardon, L. (1994). Computational DNA sequence analysis. *Annu. Rev. Microbiol.* **48**, 619–654.

Karlin, S. & Dembo, A. (1992). Limit distributions of maximal segmental score among Markov-dependent partial sums. *Adv. Appl. Prob.* **24**, 113–140.

Karlin, S., Dembo, A. & Kawabata, T. (1990). Statistical composition of high-scoring segments from molecular sequences. *Ann. Stat.* **18**, 571–581.

# References: Sequence alignments with gaps.

Gotoh, O. (1982) An improved algorithm for matching biological sequences. *J. Mol. Biol.* **162**, 705–708.

Naor, D. and D.L. Brutlag (1994) On near-optimal alignments of biological sequences. *J. Comp. Biol.* ~~4, 000–000~~. $1(4) = 1-18$.?

Needleman, S.B. and C.D. Wunsch (1970) A general method applicable to the search for similarities in the amino acid sequence of two proteins. *J. Mol. Biol.* **48**, 443–453.

Smith, T.F. and M.S. Waterman (1981) Identification of common molecular subsequences. *J. Mol. Biol.* **147**, 195–197.

Zuker, M. (1991) Suboptimal sequence alignment in molecular biology. Alignment with error analysis. *J. Mol. Biol.* **221**, 403–420.

# Genomic Signatures

Dinucleotide Relative Abundances -- invariant w/in species

$$P_{xy}^* < 0.77 \qquad P_{xy}^* : 0.97-0.81 \qquad 0.81-1.18 \qquad 1.18-$$

$$-- \qquad\qquad - \qquad\qquad 0 \qquad\qquad + \qquad\qquad ++$$

## Distances

residual dinuc $= \dfrac{f_{ij}^*}{f_i^* \cdot f_j^*} - 1$

$- \rho^*(f,g) = \sum \left| P_{IJ}^*(f) - P_{IJ}^*(g) \right| w_{ij} \quad \Rightarrow \quad w_{IJ} = \dfrac{1}{16}$

$-$ $\underline{\underline{use}}$

① Ind. Ident Distrib

$$P_{xy}^* \rightarrow 1$$

② betw. species distances generally greater than
w/in species distances

#'s v. low
w/in a species
<u>E. coli</u>

|   | 1 | 2 | 3 | 4 |
|---|-----|-----|-----|-----|
| 1 | 0 |   |   |   |
| 2 | 0.025 | 0 |   |   |
| 3 | 0.008 | 0.017 | 0 |   |
| 4 | 0.032 | 0.014 | 0.025 | 0 |
| 5 | 0.009 | 0.021 | 0.008 | 0.027 |

USES AVERAGE VALUES

ADVANTAGE OF THIS DISTANCE
① di, tetra & tri
② $\rho^* \times N_{RY}$   any $k$ this
relaxation
approaches
1

1|3|5   2|4
5|1|3

Explanations

(A) why is TA low?

    ① lowest stacking energy of dinucs
    - most unstable dinuc.

but doesn't explain non-coding bias ← ② RNase preferentially degrade UpA

    ③ part of many regulatory sequences

(B) why is CG low?

    ① methylation → deamination

$$NH_2 - CpG$$
$$\downarrow CH_3$$
$$NH_2 - CpG$$
$$\downarrow CH_3$$
$$TpG$$
$$\vdots$$
$$\downarrow$$

    - but doesn't explain why CpG low
    in mitochondria because no methylation

(C) CTAG low

    ① kinks easily
    ② clusters in rRNA
    ③ dense in replication origins

(D) GATC low in almost all bacteriophage in E.coli

If

$$\rho_{xy}^* \leq 0.78 = \text{highly significant underrepresentation}$$

$$\rho_{xy}^* \geq 1.22 = \quad '' \quad '' \quad \text{overrepresentation}$$

Examples

- CG is low for all vertebrates

- single mutants from this

$$CG \rightarrow CA \quad \Big\} \text{overrepresented in animals}$$
$$\quad \rightarrow TG$$

② Proks

- TA v. low

- C(TAG) (G̅ ̅A̅ ̅T̅.̅C̅) v. low    .... embedded in this is 2 stop codons
  ↓
- ATAG v. high    ... but other stop codons not low.

$$\frac{\overline{f_{xyz}}}{f_x f_y f_z \left(\frac{f_{xy}}{f_x f_y}\right)\left(\frac{f_{yz}}{f_y f_z}\right)\overbrace{\left(\frac{f_{x\neq z}}{f_x \cdot f_z}\right)}^{xNz}}$$

factoring
out
mononucl.
freqs

← residual dinucleotide effects

$$\frac{f_{xyz}\ f_x f_y f_z}{f_{xy}\ f_{x\cdot z}\ f_{yz}} = \text{TRINUCLEODIE FREQS w/}$$
$$\text{ALL MONO, } \underline{DI} \text{ REMOVES}$$

$$\frac{f_{xyzw}}{f_x f_y f_z f_w \left(\begin{array}{c}\text{dinucleotide}\\\text{residuals}\end{array}\right)\left(\begin{array}{c}\text{trinucleotide}\\\text{residuals}\end{array}\right)} = \text{TETRA w/ ALL MONO, DI}$$
$$\text{TRI. REMOVED}$$

REMOVING STRANDS - SYMMETRIZING

$$f_{xy}^* \qquad f_A^* = f_T^* = \frac{1}{2}\left(f_A + f_T\right)$$

. TAKE LONG
ENOUGH SEQUENCE
THEN THE
$f_A = f_T$ ,

$$f_{GT}^* = f_{AC}^* = \frac{1}{2}\left(f_{GT} + f_{AC}\right)$$

$$\therefore p_{xy}^* = \frac{f_{xy}^*}{f_x^* f_y^*}$$

# Dinucleotides

① How measure bias

    ⓐ longer is better    >19000 needed

    ⓑ <u>mononucleotide</u> content

        - GC varies between 10-80%

        - doesn't fit w/ habitat

        · must factor <u>out</u> mononucleotides

$$- E(f_{xy}) = f_x \cdot f_y$$

ⓒ

$$\boxed{\begin{array}{c} \text{-ODDS} \\ \text{RATIO} \end{array} \quad \frac{f_{xy}}{f_x f_y} - 1} = \text{DINUCLEOTIDE BIAS}$$

ⓓ <u>trinucleotides</u>

    - Markov model of order 1

$$f_{ABC} = Pr\{A|BC\} \; Pr\{BC\} \quad\quad \leftarrow \text{since order 1}$$
$$\phantom{f_{ABC} = Pr\{A|BC\}} \underset{\text{given}}{} \quad\quad C \text{ has no effect} \atop \text{on } A$$

$$= Pr\{A|B\} \; Pr\{BC\}$$

$$= \frac{Pr\{AB\} \; Pr\{BC\}}{Pr\{B\}}$$

<u>ORDER 2</u>
$$f_{ABCD} = \frac{Pr\{ABC\} \; Pr\{BCD\}}{Pr\{BC\}}$$

<u>ORDER n-2</u>
$$f_{X_1 \cdots X_n} = \frac{f\{X_1 \cdots X_{n-1}\} \; f\{X_2 \cdots X_n\}}{f\{X_2 \cdots X_{n-1}\}}$$

_Mathematical & Computational MolBio - Karlin_

<u>Enormous database of sequences</u>
  DNA - 200,000,000 bit
  PROT --- 60,000 bit
  Structures

<u>Lots</u> of opportunities to see patterns in data

  ① <u>complete genomes</u>
        - which are genes
        - <u>repeats</u>
        - <u>word patterns</u>

  ② <u>short words</u>
        - over/under representations & biases
        - <u>compare</u> diff. pieces of DNA

  ③ <u>mononucleotides</u>
        ① G+C varies enormously
             - e.g. isochores in humans
             - <u>not</u> biased in bacteria, flies

Blaisdell        ② <u>Markov models</u> --- of any order
                     - show that DNA's <u>are</u> not dependent on
                           neighboring bases
                     - too many <u>local</u> repeats

        ③ <u>Chaos models</u> --- not yet applied well

        ④ <u>Linguistic models</u> --- doesn't think this is good either
             - freq. of words

    - DOESN'T BELIEVE IN FITTING MODELS
    - BELIEVES IN BENCHMARKS

# Math./Comp. Molecular Biology

## Substitution Matrices

- substitution scores $\{S_{ij}\}$ $\quad i = 1 \cdots 20$, $j = 1 \cdots 20$ ; symmetrical
- matrices -- do <u>not</u> have to be square <u>NOR</u> symmetrical

---

<u>EXAMPLE</u> - COMPARE STRUCTURES & SEQUENCES = INVERSE FOLDING

① <u>derive counts</u>

|  | A | C | D | $\cdots$ | Y | / amino acids |
|---|---|---|---|---|---|---|
| $\alpha$ | | | | | | $\mathcal{E}_1$ |
| $\beta$ | matrix $c_{ij}$ = freq. of each a.a. in each class | | | | | $\mathcal{E}_2$ |
| $\uparrow$ | | | | | | $\mathcal{E}_3$ |

<u>structures</u> / <u>environments</u>

$f_1 \quad f_2 \quad f_3 \quad \cdots \quad f_{20}$

② form scores $\quad S_{ij} = \ln \dfrac{c_{ij}}{\Sigma_i f_i f_j} = \ln \dfrac{observed}{expected}$

③ use for comparisons

---

- $S_{ij}$ occurs w/ prob. $p_i p_j$

① in HSSPs substitution freqs are $\sim q_{ij} = p_i p_j e^{\lambda S_{ij}}$

② typical length of HSSP at a given significance $\ell$ is

$L_p = \dfrac{S_p}{H/\lambda}$

∴ as H incr. the length shrinks
as H decr. the length increases

$S_p = \begin{array}{l} significant \\ score \\ thresholds \end{array} = \dfrac{1}{\lambda} \{ \ln MN + \ln k - \ln[-\ln(1-p)] \}$

$H = \lambda \sum_{i=1}^{20} \sum_{j=1}^{20} q_{ij} S_{ij} = $ "relative entropy" (Altshul)

-H varies among matrices

③ $F_I$ = expected fraction of identities in HSSP

$$F_I = \sum_{i=1}^{i=20} g_{ii}$$

$$F_P = \sum_{i=1}^{i=20} \sum_{j=1}^{j=2} g_{ij} \quad (g > 0) \quad = \text{fraction of conservative substitutions}$$

## CREATING MATRICES

① $PIM_F$ = percent identity matrices
   = target HSSP where $F_I = F$

ⓐ start w/ subs. counts from learning set $\Rightarrow C_{ij}$
   · derive from segments w/ a lot of identities
   · assume ratio of changes will be constant even w/ lower identities

Conditionals $g_{ii} = \dfrac{C_{ii}}{\sum C_{ii}}$ $\qquad h_{ij} = \dfrac{C_{ij}}{C_{..} - \sum C_{ii}}$

$$S_{ii} = \ln \frac{F g_{ii}}{f_i \cdot f_i} \qquad\qquad S_{ij} = \frac{(1-F) h_{ij}}{f_i \cdot f_j}$$

ⓑ $\sum_{i=1}^{20} \sum_{j=1}^{20} f_i f_j \, e^{\lambda S_{ij}} = 1$

ⓒ $F_I = \sum_{i=1}^{20} g_{ii} = \sum_{i=1}^{i=20} f_i f_i e^{S_{ii}} = \sum_{i=1}^{20} F g_{ii}$

## BLAST

-also includes some correction for multiple sequences in the database

-adjusts $L$ so that length is based on length of entire databases

## PAIRWISE ALIGNMENTS

SSPA - significant segment pair alignment

① determine all HSSPs
② order HSSPs optimally
③ eliminate overlaps
④ score $= \dfrac{\Sigma HSSPs}{max. self score}$  or  $\dfrac{\Sigma}{min self score}$  ⟵ GLOBAL SCORE

$$score = \dfrac{\Sigma HSSPs}{max \; or \; min \; for \; range \; of \; alignment}$$

$$score = \dfrac{\Sigma HSSPs}{aligned \; region}$$

# Phylogeny
- taking sequences and ordering a tree
- organizing relationships
    to reflect evolutionary descent

ⓐAssume evolutionary changes are caused by mutations
    that are substitutions or dindels or inversions

## Complications
ⓐobserved now - infer past

## 3 types of similarities
- ancient shared characteristics
- derived shared characteristics
- convergent shared characteristics

## Trees

### 3 species

edges... length represents divergency/distance

$L1A + L2A$
$L1A + LA3$
$L2A - LA3$



$D_{12} \quad D_{13} \quad D_{23}$

$$L1A + L2A = D12$$
$$L1A + LA3 = D13$$
$$L2A + LA3 = D23$$

Phylogeny

$$L1A = \frac{1}{2} \{ D12 + D13 - D23 \}$$
$$L2A = \frac{1}{2} \{ D12 + D23 - D13 \}$$
$$LA3 = \frac{1}{2} \{ D13 + D23 - D12 \}$$

UNIQUE SOLUTION

-_but_ unique solution is meaningless

## 4 species



$\binom{4}{2}$ distances = **6**

$D_{12}$ $D_{13}$ $D_{14}$

$D_{23}$ $D_{24}$

$D_{34}$

## 5 Branches

$L1A$
$L2A$
$LAB$
$LB3$
$LB4$

① system is overdetermined

② but ... other trees ... 3 total    $2^{n-2}+1$



## In general

| | | $n$ | 2 | 3 | 4 | 5 | 6 | $n$ | $n+1$ |
|---|---|---|---|---|---|---|---|---|---|
| -n species | | | | | | | | | |
| -$D_n$ pairwise distances | $\binom{n}{2}$ | | 1 | 3 | 6 | 10 | 12 | $D_n$ | $D_n + n$ |
| -$L_n$ branch lengths | $2n-3$ | | 1 | $3 \cdot 2 = 5$ $\cdot$ $2 = 7$ | | | | $L_n$ | $L_n + 2$ |
| -$T_n$ topologies | $1 \cdot 3 \cdot 5 \cdots 2n-5$ | 1 | 1 | 3 | 15 | | | $T_n$ | $L_n T_n$ |

$$\frac{(2n-5)!}{(2^{n-3})(n-3)!}$$

## 4 species



- even w/ unequal rates -- can factor out rates
- because have a shared segment

$$L_{1A} + L_{A2} = D_{12}$$
$$L_{1A} + L_{AB} + L_{B3} = D_{13}$$
$$L_{1A} + L_{AB} + L_{B4} = D_{14}$$
$$L_{2A} + L_{AB} + L_{B3} = D_{23}$$
$$L_{2A} + L_{AB} + L_{B4} = D_{24}$$
$$L_{B3} + L_{B4} = D_{34}$$

what is the condition for solution to exist?

①* $D_{12} + D_{34} \leq D_{13} + D_{24} = D_{14} + D_{23}$

> $<$ $\asymp$ $\asymp$

① If $D_{13} + D_{24} = D_{14} + D_{23} \Rightarrow$ there is a unique solution with properties *

② But generally use least-squares estimate to approximate solution $L_{ij}$

i try to minimize $\underline{LS}$

## LSQ solution

$$L_{1A} = \tfrac{1}{4} \{D_{13} + D_{14} + D_{23} - D_{24}\} + \tfrac{1}{2} D_{12}$$

$$L_{2A} = \tfrac{1}{4} \{D_{23} + D_{24} - D_{13} - D_{14}\} + \tfrac{1}{2} D_{12}$$

$$L_{B3} = \tfrac{1}{4} \{D_{13} + D_{23} - D_{14} - D_{24}\} + \tfrac{1}{2} D_{34}$$

$$L_{B4} = \tfrac{1}{4} \{D_{14} + D_{24} - D_{13} - D_{23}\} + \tfrac{1}{2} D_{34}$$

$$L_{AB} = \tfrac{1}{4} \{D_{13} + D_{14} + D_{23} + D_{24}\} - \tfrac{1}{2} (D_{12} + D_{34})$$

if $D_{13} + D_{24} = D_{14} + D_{23} =$

For additive tree

$$L_{1A} = \frac{1}{2}(D_{12} - D_{23} + D_{13})$$

$\vdots$

When sequences are v. similar most mutations are <u>independent</u> ∴ distances are additive

## Fitch
- assumes relaxed additivity
- take 4 species at a time

## Methods
① Find correct tree & estimate branch lengths

distance

parsimony

likelihood

## Distance matrix methods

= n species : (OTUs operational taxonomic units)

- $\binom{n}{2}$ distances

- $\dfrac{(2n-5)!}{2^{n-3}(n-3)!}$  distinct bifurcating trees w/ OTU leaves $\qquad \left( 2^{\frac{123!}{61} \cdot 61} \right)$

- $2n-3$  branches

## Task
① find correct topology
② estimate branch lengths

## Methods
① evaluate all trees -- use LS -- get distances ---pick smallest

UPGMA = unweighted pair group method of arithmetic mean

Observed
Distances

|  | 1 | 2 | 3 | 4 |  |
|---|---|---|---|---|---|
| 1 | $d_{11}$ | $d_{12}$ | $d_{13}$ | $d_{14}$ | ... |
| 2 |  | | $d_{23}$ | $d_{24}$ | ... |
| 3 |  | | | | |
| 4 |  | | | | |

Recursion (reduce # of species

→ ① find smallest distance ... $D_{12}$

② $D_{12}/2$ ⌐⌐ $D_{12}/2$ — make branches joined w/ length $D_{12}/2$

③ replace 1 & 2 by 12

④ convert $D_{X-12} = \dfrac{D_{X_1} + D_{X_2}}{2}$

⑤

Example

|  | H | C | G | O | B |
|---|---|---|---|---|---|
| H |  | 0.094 | 0.14 | 0.180 | 0.207 |
| C |  | | 0.115 | 0.194 | 0.218 |
| G |  | | | 0.188 | 0.218 |
| O |  | | | | 0.218 |
| B |  | | | | 0.216 |

↓

|  | HC | G | O | B |
|---|---|---|---|---|
| HC |  | 0.1275 | 0.187 | 0.212 |
|  | | | 0.188 | 0.215 |
|  | | | 0.188 | 0.218 |
| G |  | | | |
| O |  | | 0.188 | 0.216 |
| B |  | | | |

→

|  | O | B |
|---|---|---|
| HCG |  | 0.188 | 0.218 |
|  | | | 0.216 |

$\frac{1}{2}\left(\begin{array}{c} HB \\ + \\ CB \\ + \\ GB \\ + \\ OB \end{array}\right)$ / 4

## ADVANTAGES
- v.v. fast

## Properties
- branch lengths always positive
- from any point internal to the leaf you get the same average distance
- ∴ only works best w/ constant evolutionary rate
- if get correct topology then the estimates of the branch lengths are least squares if assume constant rate



$1 2 > 6$
$1 3 > 0$
$2 3 > 0$

$AB > 0$

## Fitch & Margoliash

- works on triplets



$$✦ \quad L_{1A} = \tfrac{1}{2}\{D_{12} + D_{13} - D_{23}\}$$
$$L_{2A} = \tfrac{1}{2}\{D_{12} + D_{23} - D_{13}\}$$
$$L_{3A} = \tfrac{1}{2}\{D_{13} + D_{23} - D_{12}\}$$

① start w/ matrix

② find smallest Distance $D_{12}$

③ group remaining distances into $1 \rightarrow n$

$$D_{1n} = \frac{D_{13} + D_{14} + D_{15} \cdots D_{1n}}{n-2}$$

$$D_{2n} = \frac{D_{23} + D_{24} + D_{25} \cdots D_{2n}}{n-2}$$

④ calculate branch lengths from ✦

⑤ group 1, 2

⑥ then convert $D_{1n}$ $D_{2n}$ like UPGMA ···w/ averages

## Example --

$$D_{HC} = 94$$
$$D_H(GOB) = \tfrac{1}{3}(111 + 180 + 207) = 166$$
$$D_C(GOB) = 176$$



$\tfrac{1}{2}(D_{12} + D_{23} - D_{13}) = 52$



$$\bigstar = 55 - \frac{(42+2)}{2}$$

rooting



Assume constant rate

$$2\lambda t = r$$
$$= x + z - r$$
$$= y + z - r$$

Least Squares $\quad \bar{r} = \frac{1}{4}(x + y + 2z)$

Then -- do branch swapping
    ·· try to minimize Least Squares

<u>Neighbor - Joining</u>   Saitou & Nei 1987 JME 4:406-425

① n species w/ matrix
② derive topology differently
③ begin w/ star-like topology



$= L_{1x}, L_{2x}, L_{3x}, L_{4x} \cdots L_{nx} =$ branches

$=$ distances $= D_{ij}$

④ want to minimize sum of branch lengths
⑤ assume branch lengths are additive

$$S_0 = \sum_{i=1}^{n} L_{xi} \overset{\text{ASSUME}}{\underset{\text{ADDITIVE}}{=}} \frac{1}{n-1} \sum_{i<j} D_{ij}$$



⑥ consider diff. topologies



length $= L_{12} \quad L_{1x} \quad L_{2x} \quad L_{xy} \quad L_{y3} \cdot L_{y4}$

SUM FOR  = $S_{12} = \underbrace{L_{1x} + L_{2x}}_{D_{12}} + L_{xy} + \underbrace{\sum_{i=3}^{n} L_{yi}}_{\frac{1}{n-3}\sum D_{ij}} \overset{\text{ASSUME}}{\underset{\text{ADDITIVE}}{=}} D_{12} + \frac{1}{n-3} \sum_{3 \le i \le j}^{n} D_{ij}$
THIS TOPOLOGY

$+ \frac{1}{2(n-2)} \sum_{k=3}^{n} (D_{1k} - L_{1x} - L_{yk} + D_{2k} - L_{2x}$
$\qquad\qquad - L_{y2})$

$$= \frac{1}{2(n-2)} \sum_{k=3}^{n} \left( D_{1k} + D_{2k} \right) + \frac{1}{2} D_{12} + \frac{1}{n-2} \sum_{3 \leq i < j} D_{ij}$$

ⓔ calculate <u>all</u> possible $D_{xy}$ $\binom{n}{2}$

ⓔ choose that which gives minimum sum

ⓕ reduce distance matrix by taking averages

It turns out — that the distances come out to be
the least squares estimate.

## Neighbor - Joining



$S_{12} = $ sum of branch lengths

$$= L_{1x} + L_{2x} + L_{xy} + \frac{1}{2}\sum_{k=3}^{n} L_{yk}$$

$$= \frac{1}{2}D_{12} + \frac{1}{2(n-2)}\sum_{i=3}^{n}\left(D_{1i} + D_{2i}\right) + \frac{1}{n-2}\sum_{3 \le i < j} D_{ij}$$

$$Min\left[\left(L_{12} - D_{2}\right)^2 + \left(L_{13} - D_{3}\right)^2 + \left(L_{14} - D_{4}\right)^2 \cdots\right]$$

$$\sum \hat{L}_{1x} = \sum \hat{S} = S_{12}$$

Calculates branch lengths as with Fitch-Margoliash

## Claim

NJ---if distances are additive it will find correct topology.

(A)

$$S_{12} = \frac{1}{2} D_{12} + \frac{1}{4}\left(D_{13} + D_{14} + D_{23} + D_{24}\right) + \frac{1}{2} D_{34}$$
$$= S_{34}$$

$$S_{13} = S_{24} = \frac{1}{2} D_{13} + \frac{1}{2} D_{24} + \frac{1}{4}\left(D_{12} + D_{14} + D_{23} + D_{34}\right)$$

$$S_{13} - S_{12} = -\frac{1}{4} D_{12} + \frac{1}{4} D_{13} + \frac{1}{4} D_{24} - \frac{1}{4} D_{34}$$
$$= \frac{1}{4}\left(D_{13} + D_{24} - D_{12} - D_{34}\right)$$

$S_{13} - S_{12} > 0 \Rightarrow$ tree A is correct

$D_{12} + D_{34} < D_{13} + D_{24}$        u.similar to before for additive tree
$$D_{12} + D_{34} \leq D_{13} + D_{24} = D_{14} + D_{23}$$

∴ for $N=4$  NJ gives correct tree

u/ non-additive

$$D_{12} + D_{34} \leq D_{13} + D_{24} \leq D_{14} + D_{23}$$
$$D_{12} + D_{34} \leq D_{14} + D_{23} \leq D_{13} + D_{24}$$

one will be tree
if

NJ
.1 species
_____

.induction from <u>n-1</u> to n
-3 possibilities

$$2\,'\!\!\!\underset{}{\diagdown}\!\!\!\overset{}{\underset{}{\searrow}}\!\!\!\text{⋯} \quad\text{or}\quad \underset{new}{\overset{1}{\diagdown}}\!\!\!\overset{}{\underset{}{\searrow}}\!\!\!\text{⋯} \quad\text{or}\quad \underset{new}{\overset{2}{\diagdown}}\!\!\!\overset{}{\underset{}{\searrow}}\!\!\!\text{⋯}$$

— If $S_{12}$ is the smallest
consider $S_{in} - S_{12} = \frac{1}{2}D_{1n} + \frac{1}{2(n-2)}\Big[\big(D_{12} + D_{13} + \cdots D_{1n-1}\big) + \big(D_{2n} + D_{3n} + \cdots D_{3n_1}\big)\Big]$

$$+ \frac{1}{n-2}\Big[\big(D_{23} + D_{24} \cdots D_{3,n-1}\big) + \big(D_{34} + \cdots D_{3,n-1}\big) + \big(\cdots \big)\Big]\Big]$$

$$- \frac{1}{2}D_{12} - \frac{1}{2(n-2)}\Big[\big(D_{13} + D_{14} \cdots D_{1n}\big) + \big(D_{23} \cdots D_{2n}\big)\Big]\Big]$$

$$+ \frac{1}{n-2}\Big(D_{34} + D_{35} + \cdots D_{3n\,-1} + D_{45} + \cdots D_{4n}\Big)$$

$$\boxed{S_{in} - S_{12} = \frac{1}{2n-2}\sum_{k=3}^{n-1}\big(D_{1n} + D_{2k} - D_{12} - D_{kn}\big)}$$

~~$S_{in} - S_{12}$~~
<u>Assume</u>  $\underset{n}{\overset{1}{\diagdown}}\!\!\!\!\overset{a}{\underset{b}{\diagup}}\!\!c\!\!\overset{d}{\underset{e}{\diagdown}}\!\!\!\!\overset{2}{\diagup}_{k}$

$$D_{1n} + D_{2k} - D_{12} - D_{kn} = a+b+d+e - (a+c+d) - (b+c+e)$$
$$= -2c$$

$$\therefore \sum_{k=3}^{n-1} -2c \quad\text{is always} < 0$$

$$\therefore S_{in} - S_{12} \text{ would be } > 0$$

Fitch, Doolittle & Feng          (JME 18:30.1981)

Neighborliness

- give two neighbors a $\pm 1$ score for __each__
  set of four sequences with that pair
  in which those __two__ are closest

  ~~##~~ $\binom{n-2}{2}$  # of choices

- use this matrix of neighborliness for tree making

<u>Parsimony</u>

① Given an alignment --- each aligned position is represented by a tree
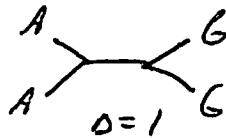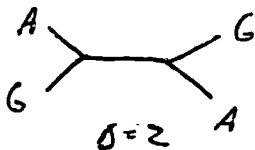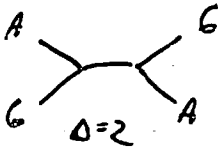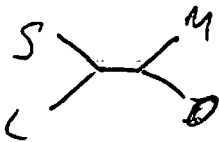② Reconstruct internal nodes of tree & get fewest substitutions
③ Add up over all columns

Predict True Tree = tree w/ minimal substitutions

<u>But</u> --- must look at all trees

<u>EXAMPLE</u>

| | Non Inf | Non Inform. | | | | Non Inform | | Non Inform | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| shark | G | A | T | C | C | T | A | G | G | C |
| lungfish | G | G | T | C | A | C | A | T | G | T |
| monkey | G | G | T | C | A | T | A | T | C | T |
| outgroup | G | A | T | A | C | C | A | G | C | A |
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |



$\Sigma = 9$      $\Sigma = 9$      $\Sigma = 7$

Maximum likelihood

Phylogenetic _or_ phenetic trees

Principle ...

Observed data - D

- several alternative probability models (e.g. ... diff. trees) $M_i$
- $P(D|M_i)$ = prob. of observing $D$ under model $M_i$
- $P(M_i|D)$ = likelihood of $\underline{M_i}$ given

$$= \frac{P(M_i, D)}{P(D)} = p \, \frac{P(D|M_i) \, P(M_i)}{\sum\limits_{i} P(D|M_i) P(M_i)} = \begin{array}{l}\text{BAYES}\\ \overline{\text{FORMULA}}\end{array}$$

$P(M_i)$ = a priori probabilities

$\therefore$ assume all $P(M_i)$ equally likely

$$= \frac{P(D|M_i)}{\sum P(D|M_i)}$$

$\therefore$ Most likely model is the one that maximizes $P(D|M_i)$

Example

- coin tossing
  - if $prob(H) = p = M_p$
  - f $n$ tosses w/ $k$ heads = $D$
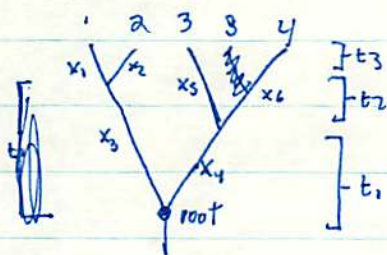
$$L = Prob(D|M_p) = \binom{n}{k} p^k (1-p)^{n-k}$$

$$\Downarrow$$

$$\log L = \log\binom{n}{k} + \log p^k + \log(1-p)^{n-k} \qquad \text{use } \underline{\underline{\log}} \text{ Likelihood}$$

$\underset{\text{WHEN}}{\lambda, \beta'} = \dfrac{d \log L}{d\rho} = .0. = \cdots$

$$\hat{p} = \dfrac{k}{n}$$

## MAX. LIKELIHOOD IN PHYLOGENY

① GIVEN TREE ⋯ ESTIMATE BRANCH LENGTHS

(LANGLEY & FITCH JME 3:161)



- $x_i = $ # of substitutions
- times = unknown

model = poisson process for substitution events

$$P_t(X = x) = \text{prob. of } x \text{ subs. in time } t$$
$$= \dfrac{\lambda t^x e^{-\lambda t}}{x!}$$

### ASSUMPTIONS
- events in one time period independent of others
- linearity for small times
- $\lambda$ is constant

### ESTIMATION
- can only calculate $\lambda^t = v$

LIKELIHOOD $= \prod L_1 \cdot L_2 \cdot L_3 \cdot L_4 \qquad L_1 = \dfrac{v_1^{x_4} e^{-v_1}}{x_4!} \qquad L_2 = \dfrac{(v_1 + v_2)^{x_3} e^{-(v_1 + v_2)}}{x_3!}$
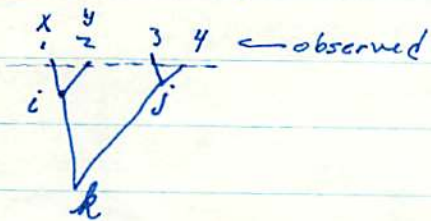
$$L_3 = \dfrac{v_3^{x_1 + x_2} e^{-2v_3}}{x_1! \, x_2!} \qquad L_4 = \dfrac{(v_2 + v_3)^{x_5 + x_6} e^{-2(v_2 + v_3)}}{x_5! \, x_6!}$$

MAX LIKELIHOOD = take derivative $\dfrac{d}{dx_1}, \dfrac{d}{dx_2}, \dfrac{d}{dx_3}, \dfrac{d}{dx_4}$

= set to zero

= solve

2) GIVEN SEQS.-- ESTIMATE TOPOLOGY & BRANCH LENGTHS

- specify ... $P_{ij}^t$ = prob. in time $t$ that get $i \to j$ change

- likelihood can be specified for <u>any</u> topology ... on a per nucleotide basis



1  2  3  4  ← observed

$$L = \sum_{k=1}^{4} \left( P_k \cdot \left( \sum_{i=1}^{4} P_{ki}^{t_i} P_{ix}^{t_d} P_{iy}^{t_2} \right) \left( \sum_{j=1}^{4} P_{kj}^{t_3} P_{jz}^{t_d} P_{jv}^{t_4} \right) \right)$$

<u>One parameter Model</u>   $P_{ij}^t = \left(1 - e^{-\lambda t}\right) P_j$    $i \neq j$

$P_{ii}^t = e^{-\lambda t} + \left(1 - e^{-\lambda t}\right) P_i$           $e^{-\lambda t}$ = prob of <u>no</u> event

<u>Algorithm</u>
   ① sum up log-likelihoods over all sites
   ② find maximum in terms of $v = \lambda t$
   ③ get most likely tree

<u>Consistency checks</u>
   ① bootstrapping ... resample columns w/ replacement
       - but...

## Pattern freq distribution

$$P\left(\begin{array}{c}\text{pattern } P \\ \text{occurs } x \text{ times}\end{array}\right) = ?$$

Definition: Sequence $S$ is in state $(x, i)$ if it contains $P$ $x$ times and ends on $P_i$

$$i = 0 \dots m-1$$

## State $\Delta$ $f(x)$

$$s(x, i : j) = \begin{cases} x_i, & t_{ij} \\ x+1, & i_x \end{cases} \begin{array}{l} t_{ij} \neq m \\ t_{ij} = m \end{array}$$

To study words depends on what order of markov model
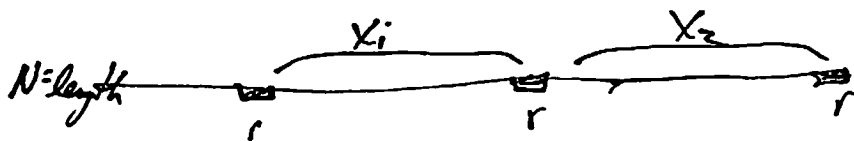
### 1 ORDER

- word = $AAAT$
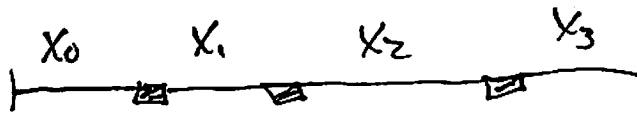- need A
      C
      G
      T
      AA
      AAT

## Distributions between patterns



$N = $ length

$K = $ # of events

$$P(N_m < k) = P(X_1 + \dots X_{k-1} \geqslant m)$$

- if $N$ is large, $k$ large   $X = $ normally distributed

$$X_0 \quad X_1 \quad X_2 \quad X_3$$

$X_0 =$ 1st passage

How get Markov model

# Multiple Alignments

1) Lawrence & Reilly Proteins.
7:41

2) Krogh et al JMB 235:1501

3) Lawrence et al Science 263:208

## global multiple alignments
- progressive pairwise alignments
- profile methods

## local multiple alignments = motifs??
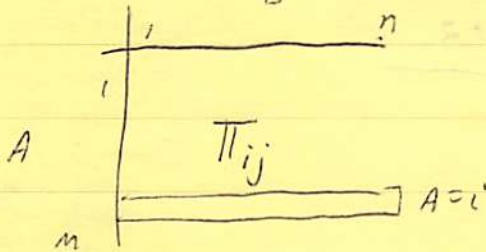
e.g - Protein - DNA binding sites
- promoters

Cannot do pairwise alignments because each p.w. will come up w/ diff. regions

## TWO FUNDAMENTAL CONCEPTS

① conditional probability
- two events $A$ & $B$



- cond. prob = prob $\{B=j \mid A=i\}$

$= prob \{B=j, A=i\}$

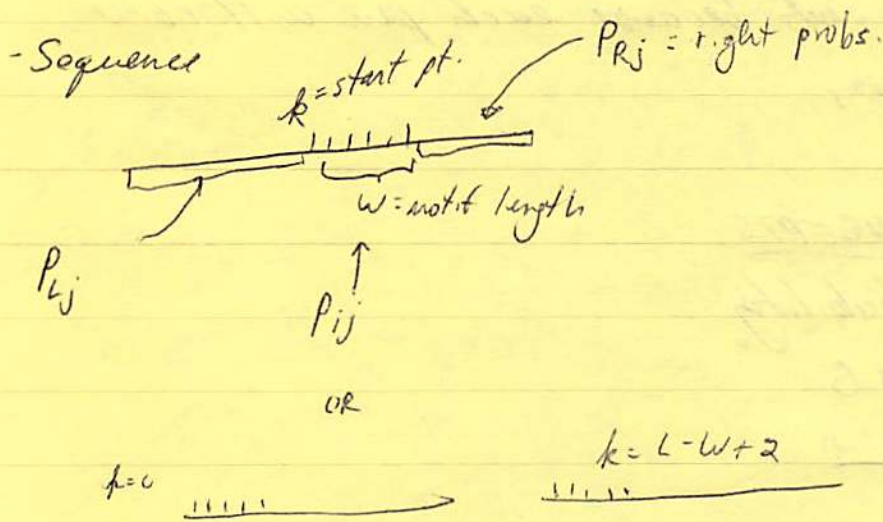$= \dfrac{\pi_{ij}}{\sum\limits_{j=1\cdots n} \pi_{ij}}$

$$\text{(2)} \quad f(x_1, x_2 \ldots x_n) = \sum_{i=1}^{i=n} p_i \log x_i \quad w/ \quad \sum x_i = 1$$

this $f(x)$ has maximum $\quad x_i = \dfrac{p_i}{\sum\limits_{i=1}^{n} p_i}$

A) - Motif has no gaps
 - Some seqs have motif; some don't

$\qquad$ n seqs $\{S_1 \cdots S_n\}$
$\qquad$ length $\{L_1 \cdots L_n\}$

 - Sequence $\qquad k = \text{start pt.} \qquad P_{Rj} = \text{right probs.}$



$\qquad\qquad\qquad W = \text{motif length}$

$P_{Lj} \qquad\qquad\qquad P_{ij}$

$\qquad\qquad\qquad$ OR

$k = 0 \qquad\qquad\qquad\qquad k = L - W + 2$

$P_\sigma = $ prob. that seq has motif

want to maximize $\quad \overset{\text{likelihood}}{\lambda_\theta(s)} \qquad \theta = P_J, P_{ij}, P_{Rj}, P_\sigma$

## EM Algorithm

- start w/ $\theta$
- finds $\theta'$ such that $g(\theta')$ is $\geq g(\theta)$
- the way to find the next $\theta'$

$$\log g_\theta(s_s) = \log h_\theta(s_s, k) - \log w_\theta(k \mid s_s)$$

$$\log g_{\theta'}(s_s) - \log g_\theta(s_s) = \frac{\log h_{\theta'}(s_s, k)}{h_\theta(s_s, k)} - \log \frac{w_{\theta'}(k \mid s)}{w_\theta(k \mid s)}$$

- multiply by $w_\theta(k \mid s_s)$ & sum ($=1$ because $\Sigma = 1$)

$$\log g_{\theta'}(s_s) - \log g_\theta(s_s) = \sum_k w_\theta(k \mid s_s) \log \frac{h_{\theta'}(s_s, k)}{h_\theta(s_s, k)}$$

$$- \underbrace{\sum_k w_\theta(k \mid s_s) \log \frac{w_{\theta'}(k \mid s)}{w_\theta(k \mid s)}}_{>0 \text{ by eqn } 1}$$

$$g_{\theta'}(s_s) \geq g_\theta(s_s)$$

if $\theta'$ such that $\max \sum w_\theta(k \mid s_s) \log h_{\theta'}(s_s, k)$

$\underbrace{\phantom{xxxxxxxxxxx}}$ max from $1$ cgn

$$= \sum_j \sum_k \underbrace{w_\theta(k \mid s_s) \, n(\iota s k j) \log P'_{\iota j}}$$

$$+ (\text{for right}) + \sum w_\theta(k \mid s_s) \log(w_{s_\theta} \mid k)$$

max when $P'_{\iota j} = \dfrac{\sum_k w_\theta(k \mid s_s) \, n(\iota s k j)}{\sum_j \sum_k w_\theta(k \mid s_s) \, n(\iota s k j)}$

Algorithm for searching for max $\theta$

$$\rho_\theta(S_S \mid K) = \prod_j P_{LJ}^{n(LSkj)} \, P_{RJ}^{(nRskj)} \prod_j P_{ij}^{(nskij)}$$

over: "# of times J in left"   "is $J$ at position"

$$\rho_\theta(S_S, k) = \rho_\theta(S_S \mid K) \cdot (\text{prob mot.f is at } K)$$

$$= \rho_\theta(S_S \mid k) \, W_{S\theta}(k)$$

$$W_{S\theta}(k) = P_f\left(\frac{1}{L_S - w + 1}\right) : k = 1, 2 \cdots$$

$$\frac{1 - P_J}{2} \quad = k = 0 \;\; \underline{OR} \;\; k = w - w 2$$

$$\boxed{g_\theta(S_J) = \sum_k \rho_\theta(S_{J\theta}, k)} - MAXIMIZE$$

$$W_\theta(k \mid S) = \frac{\rho_\theta(S_{g}, k)}{g_\theta(S_J)}$$

to do this... start w/ $\theta$, search for local maxima
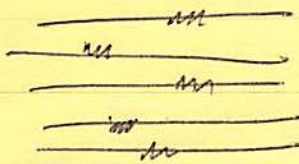
repeat

# EM Continued

- for S sequences replace $\sum_k$ w/ $\sum_S \sum_k$

- interpretation
  - for any parameters···
    - can calc. prob. that motif is at $k$ / $we(k|S_s)$
    - can then calc. new parameters from observed freqs of $i, j, k$···

## Stochastic Analog = GIBBS Sampling

try to maximize signal vs. background score

signal = $P_{ij}$
background = $P_j$

- for position $k$  $A_k = \dfrac{\pi P_{ij}}{\pi P_j}$ = background score $\quad$ = signal score

∴ want to maximize: $\sum_S \sum_i \sum_j n\, k_{s\,ij} \log \dfrac{P_{ij}}{P_j}$ over all $k$s

- ① choose random start
- ② leave 1 sequence out
- ③ scan that sequence for new $k$ using weights from other seqs
- ④ new parameters

# Doug Brutlag    Correlations → Structure in Biological Sequences

Molecular Biology is a Information Science

① $DNA → RNA → PROTEIN → FUNCTION$

② Genetic info → molecular structure → biochemical function → biological behavior

## Problems
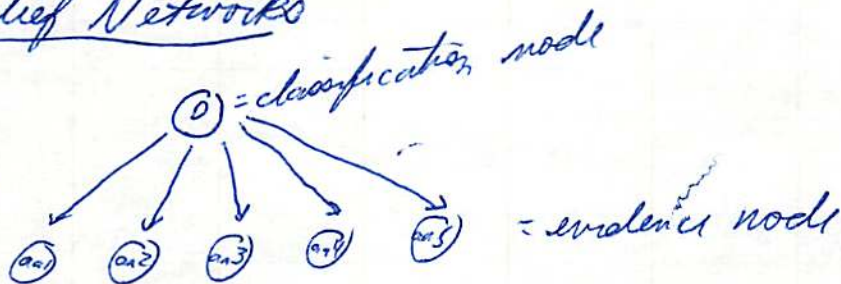- genetic info is redundant
- structural info is redundant
- multiple features encoded by 1 sequence
    - protein sequency
    - folding
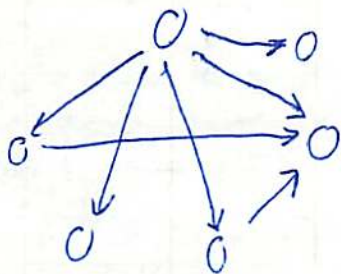    - tl & tx rate

## Representation
- most representations assume that sites are independent

## Belief Networks

- see Neapolitan



state likelihoods

= classification node

= evidence node

see Protein Science
send email
to brutlag@cmgm

can add correlation

## α-helix

- in 3D space residue $i$ is closest to $i+3, i+4$
- took a reduced set w/ ~~less~~
  - no homologs

## Test of correlations

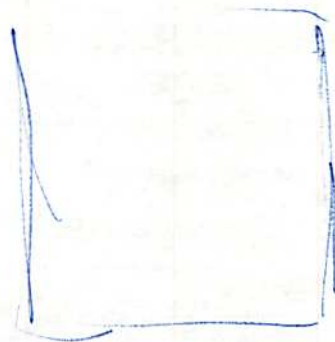① chi-squared

② mutual information

③ monte-carlo simulations

## Examples

| $i, i+4$ | D | not D |
|----------|---|-------|
| K        |   |       |
| not K    |   |       |

overrepresented          underreps

  KD   EK   SA       KL
  KĒ   FM   GA
  LL   ĪL   PF

- removed these helices
- appears these aa's interact
  - ① RANDOMIZED BONDS ···

<u>Generalized</u>

- reduced alphabet size
  by classifying into different alphabets

- convert aa into # (parametric) & look
  for correlation coefficients

(1) Do w/ 20 × 20 alphabet

repeat
- pick most similar
- group them
  (but what about forcing groups)



this is when you
want to cut off

information

- real

- random