Seq.Align. - Protein Function



Sequence Alignment

Motivation:

- •Storing, retrieving and comparing DNA sequences in Databases.
- •Comparing two or more sequences for similarities.
- •Searching databases for related sequences and subsequences.
- •Exploring frequently occurring patterns of nucleotides.
- •Finding informative elements in protein and DNA sequences.
- •Various experimental applications (reconstruction of DNA, etc.)

Exact Pattern Matching

- Given a pattern P of length m and a (longer) string T of length n, find all the occurrences of P in T.
- Naïve algorithm: O(m*n)
- Boyer-Moore, Knuth-Pratt-Morris: O(n+m)

Alignment - inexact matching

Substitution - replacing a sequence base by another.
 Insertion - an insertion of a base (letter) or several bases to the sequence.

Deletion - deleting a base (or more) from the sequence.

(*Insertion* and *deletion* are the reverse of one another)

A possible alignment of the insulin proteins from sheep and zebrafish is

Fish: MAVWLQAGALLVLLVV-SSVSTNPGTPQHLCGSHLVDALYLVCGPTGFFYNPK--R Sheep: MALWTRLVPLLALLALWAPAPAHAFVNQHLCGSHLVEALYLVCGERGFFYTPKARR

Fish: DVE-PLLGFLPPKSAQETEVADFAFKDHAELIRKRGIVEQCCHKPCSIFELQNYCN Sheep: EVEGPQVGAL--ELAGGPG-AG-GL-EGPP-Q-KRGIVEQCCAGVCSLYQLENYCN

Seq. Align. Score

SEQ 1GTAGTACAGCT-CAGTTGGGATCACAGGCTTCT||||||||||||||||SEQ 2GTAGAACGGCTTCAGTTG---TCACAGCGTTC-

SEQ 2 GINGANOGGOIICNGIIG ICNONGOGIIC

Distance 1 - match 0, substitution 1, indel $2 \Rightarrow$ distance = 14.

Distance 2 - match 0, d(A,T)=d(G,C)=1, d(A,G)=1.5 indel 2 \Rightarrow distance = 14.5.

Similarity - match 1, substitution 0, indel $-1.5 \Rightarrow$ similarity = 16.5.

General setup - substitution matrix S(i,j), indel S(i,-) or S(-,j).

Commonly used matrices: **PAM250**, **BLOSUM64**

Global Alignment

Global Alignment

INPUT: Two sequences S and T of roughly the same length. **QUESTION:** What is the maximum similarity between them? Find one of the best alignments.

The IDEA

s[1...n] t[1...m]

To align s[1...i] with t[1...j] we have three choices: * align s[1...i-1] with t[1...j-1] and *match* s[i] with t[j] s[1...i-1] i t[1...j-1] j * align s[1...i] with t[1...j-1] and *match* a space with t[j] s[1...i] t[1...j-1] i * align s[1...i-1] with t[1...j] and *match* s[i] with a space s[1...i-1] i t[1...i] -

Recursive Relation

Define: scoring matrix $m(a,b) a, b \in \sum U \{-\}$

Define: H_{ii} the best score of alignment between s[1...i] and t[1...j]

for 1 <= i <= n, 1 <= j <= m

$$H_{ij} = \max \begin{cases} H_{i-1j-1} + m(s_i,t_j) \\ H_{ij-1} + m(-,t_j) \\ H_{i-1j} + m(s_i,-) \end{cases}$$

Optimal alignment score = H_{nm}

 $H_{i0} = \sum_{0,k} m(s_{k},-)$ $H_{0j} = \sum_{0,k} m(-,t_{k})$

Needleman-Wunsch 1970

			Р	R	Е	Т	Е	R	I	Т
	i\j	0	1	2	3	4	5	6	7	8
	0	0	-1	-2	-3	-4	-5	-6	-7	-8
v	1	-1								
Е	2	-2			H _{i-1,j-1}	H _{i-1,j}				
I	3	-3			H _{i,j-1}	[♥] H _{i,j} ♥ ∽				
Т	4	-4								
G	5	-5								
Е	6	-6								
I	7	-7								
s	8	-8								
Т	9	-9								

t







Figure 1.2 The dynamic programming matrix for the example sequences, and how the values of the cells are calculated. Row and column 0 are initialized for the score of a blank equal to -1. To calculate the value of $H_{i,j}$ one needs the values of $H_{i,j-1}$, $H_{i-1,j}$ and $H_{i-1,j-1}$.

S

$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	iij 0 1 2 3 4 5 6 7 0 0 -1 -2 -3 -4 -5 -6 -7 V 1 -1 0 1 2 -3 -4 -5 -6 -7 V 1 -1 0 1 2 -3 -4 -5 -6 -7 E 2 -2 -1 0 0 -1 -2 -3 -4 I 3 -3 -2 -1 0 0 -1 -2 -2 T 4 -4 -3 -2 -1 1 0 -1 -2 -2	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	P R E T E R
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	P R E T E R I T iij 0 1 2 3 4 5 6 7 8 0 0 -1 -2 -3 -4 -5 -6 -7 -8 1 -1 0 1 2 3 4 5 6 7 8 2 -2 -1 0 0 -1 -2 -3 -4 -5 2 -2 -1 0 0 -1 -2 -3 -4 -5 3 -3 -2 -1 0 0 -1 -2 -2 -3 4 -4 -3 -2 -1 1 0 -1 -2 -1 5 -5 -4 -3 -2 0 1 -0 -1 -2 -1 5 -5 -4 -3 -2 0 1 0 -1 -2 7 -7 -6 -4	iij 0 1 2 3 4 5 6 7 0 0 -1 -2 -3 -4 -5 -6 -7 1 -1 0 1 2 3 -4 -5 -6 -7 2 -2 -1 0 0 -1 -2 -3 -4 3 -3 -2 -1 0 0 -1 -2 -2 4 -4 -3 -2 -1 1 0 -1 -2 -2	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	P R E T E R
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	P R E T E R
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	P R E T E R
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	R E T E R I T 2 3 4 5 6 7 8 -2 -3 -4 -5 -6 -7 -8 1 2 3 4 4 5 6 7 0 0 -1 -2 -3 -4 -5 -1 2 -3 4 4 5 6 7 0 0 -1 -2 -3 -4 -5 -2 -1 1 0 -1 -2 -3 -2 -1 1 0 -1 -2 -3 -3 -2 0 1 -2 -1 -2 -4 -2 -1 1 1 0 -1 -2 -5 -3 -2 0 1 2 1 2 -6 -4 -3 -1 0 1 2 1	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	R E T E R
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	E T E R
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	T E R
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	-5 -6 $-74 -2 -3 -2 -3 -2 -3 -3 -2 -3 -3 -3 -3 -3 -3 -3 -3$	E R
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$ \begin{array}{cccccccccccccccccccccccccccccccccccc$		R
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	-2 -2 -2 -2 -2 -2 -2 -2		
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	2 1 -1 -2 -3 -5 7 -8 8 T			-
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$			-5 -7 -8	Ч
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$				
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	×	<u>н</u> н ш <		
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	8 1 6 5 4 3 2 1 0 jjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjj	4 3 2 1 0 iÿ	2 2 1 0	
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$			
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	P -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1	-3 -1 -1 -1 -1 -1 -1 -1 -1		Р
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	R R -5 -5 -5 -5 -5 -5 -5 -5 -5 -5 -5 -5 -5	-1	· 0 -1 -2	R
$\begin{array}{cccccccccccccccccccccccccccccccccccc$	E -2 -2 -2 -2 -2 -2 -2 -2 -2 -2	$-\frac{3}{-3}$	-2 -3	щ
$\begin{array}{cccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$		<u>-</u>	T
$\begin{array}{cccccccccccccccccccccccccccccccccccc$	<u>-</u> 0 <u>-</u> <u>-</u> 0 <u>-</u>	0 -1 -2 -5 5	2 -4 -5	ш
	0 1	- <u>-</u> 2 - <u>-</u> 5 - <u>6</u> 6		R
	$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$		4 6 -1	Ι
				н

alignments with highest score. maximum score. Note that not all arrows are drawn. (b) The arrows showing the paths giving Figure 1.3 (a) The filled-in matrix. The arrows show which cells are used for finding the

(a)

Local Alignment

Local Alignment

INPUT: Two sequences S and T.
QUESTION: What is the maximum similarity between a subsequence of S and a subsequence of T?
Find most similar subsequences.

S = g g t c t g a g T = a a a c g aediting operations values: match = 2 indel/substitution = -1 The best *local alignment* is: $\alpha = \text{ctga} \quad (\in S)$ $\beta = \text{c-ga} \quad (\in T)$

Recursive Relation

for 1 <= i <= n, 1 <= j <= m

$$H_{ij} = \max \begin{cases} H_{i-1j-1} + m(s_i,t_j) \\ H_{ij-1} + m(-,t_j) \\ H_{i-1j} + m(s_i,-) \\ 0 \end{cases}$$

 $H_{i0} = 0$ $H_{0i} = 0$

Optimal alignment score = max_{ii} H_{ii}

Smith-Waterman 1981

Penalties should be negative

Sequence Alignment

Complexity:

Time *O*(*n***m*)

Space O(n*m) (exist algorithm with O(min(n,m)))

Ends free alignment

Ends free alignment

INPUT: Two equences S and T (possibly of different length). **QUESTION:** Find one of the best alignments between subsequences of S and T when at least one of these subsequences is a prefix of the original sequence and one (not necessarily the other) is a suffix.



Gap Alignment

Definition: A gap is the maximal contiguous run of spaces in a single sequence within a given alignment. The length of a gap is the number of indel operations on it. A gap penalty function is a function that measure the cost of a gap as a (nonlinear) function of its length.

Gap penalty

INPUT: Two sequences S and T (possibly of different length). QUESTION: Find one of the best alignments between the two sequences using the gap penalty function.

Affine Gap:

$$W_{total} = W_g + qW_s$$

 W_{g} – weight to open the gap

 $\rm W_{s}\xspace$ – weight to extend the gap

What Kind of Alignment to Use?

- The same protein from the different organisms.
- Two different proteins sharing the same function.
- Protein domain against a database of complete proteins.
- Protein against a database of small patterns (functional units)

 <u>Task</u>: Given a query sequence and millions of database records, find the optimal alignment between the query and a record

ACTTTTGGTGACTGTAC



• <u>**Tool:</u>** Given two sequences, there exists an algorithm to find the best alignment.</u>

 Naïve Solution: Apply algorithm to each of the records, one by one

• **Problem:** An *exact algorithm* is too slow to run millions of times (even linear time algorithm will run slowly on a huge DB)

• Solution:

- Run in parallel (expensive).
- Use a fast (*heuristic*) method to discard irrelevant records. Then apply the *exact algorithm* to the remaining few.

General Strategy of Heuristic Algorithms:

-Homologous sequences are expected to contain un-gapped (at least) short segments (probably with substitutions, but without ins/dels)

-Preprocess DB into some fast access data structure of short segments.

FASTA Idea

- Idea: a good alignment probably matches some identical 'words' (*ktups*)
- Example:

Database record:

ACTTGTAGATACAAAATGTG

Aligned query sequence:

A-TTGTCG-TACAA-ATCTGT

Matching words of size 4

Dictionaries of Words

ACTTGTAGATAC Is translated to the dictionary: ACTT, CTTG, TTGT,

TGTA...

Dictionaries of well aligned sequences share words.

FASTA Stage I

- Prepare dictionary for db sequence (in advance)
- Upon query:
 - Prepare dictionary for query sequence
 - For each DB record:
 - Find matching words
 - Search for long diagonal runs of matching words
 Position in
 - Init-1 score: longest run
 - Discard record if low score

*****= matching word

DB record



FASTA stage II

- Good alignment path through many runs, with short connections
- Assign weights to runs(+) and connections(-)
- Find a path of max weight
- Init-n score total path weight
- Discard record if low score



FASTA Stage III

- Improve Init-1. Apply an exact algorithm around Init-1 diagonal within a given width band.
- Init-1 → Opt-score new weight
- Discard record if low score



FASTA final stage

 Apply an exact algorithm to surviving records, computing the final alignment score.

P-Value

The observed number of random records achieving E-value *E* or better (smaller) is distributed Poisson(*E*)

Prob(*r* such records) =
$$\frac{\exp(-E)E^{r}}{r!}$$

Note: The model assumes an I.I.D. trial for each database record

BLAST (Basic Local Alignment Search Tool) Approximate Matches

BLAST:

Words are allowed to contain inexact matching.

Example:

In the polypeptide sequence IHAVEADREAM The 4-long word HAVE starting at position 2 may match

HAVE, RAVE, HIVE, HALE, ...

Approximate Matches

For each word from DB generate similar words (according to the substitution matrix) and store them in a look-up table.

BLAST Stage I

- Find approximately matching word pairs
- Extend word pairs as much as possible, i.e., as long as the total weight increases
- Result: High-scoring Segment Pairs (HSPs)

THEFIRSTL INIHAVEADREAMESIRPATRICKREAD INVIEIAMDEADMEATTNAMHEWASNINETEEN

BLAST Stage II

• Try to connect HSPs by aligning the sequences in between them:



PAM250

R Ν D С 0 Ε G Н Ι L Κ Μ F Ρ s Y V в Ζ Х * A т Ы 2 -20 0 -2 0 1 -2 0 0 0 -8 A 0 -1 -1 -1 -1 -3 1 1 1 6 з 0 R -2 6 0 -4 1 -1 -32 -3 3 0 -4 0 0 -1 2 -4 -2 -1 0 -1 -8 Ν 0 0 2 2 -4 1 1 0 2 2 -31 -2 -30 1 О 2 2 1 0 -8 2 4 2 3 1 Ο. -3Ο -2 3 3 -1 -8 D 0 -1-5 1 2 -4 -6 -1 О -5 -3 -3 0 -2 -4 -5 -3 -8 12 -5 -5 -5 -3 2 -8 0 С -6 2 Q 0 2 -54 з -21 0 -1 -5 -4 -2 1 з. -8 1 1 2 -5 - 1 Ε 0 3 2 4 0 0 0 -2 3 3 -8 -11 -5 1 2 -32 -5 -1 О 4 -1 G -3 -3 0 5 2 -2 -3 0 0 0 -1 -8 1 0 1 -3 -4 -5 1 Ο. -5 -1 -1-7 Η 2 2 1 -3 3 1 6 -2 -2 0 -2 -2 0 0 -2 1 2 -8 -1з -1 5 2 -2 2 -2 -8 Т -2-2 -3 1 2 О 5 4 L -3 -3-3 -2 2 6 -3 4 2 -3-3 2 -3 з -8 6 -2 к 3 n -5 0 -2 0 -3 5 0 -5 0 0 3 n -8 -1 1 1 2 4 0 М 0 3 -5 -3 6 0 4 2 2 -8 2 -5 0 9 -5 7 F -3-3 -5 5 -5 -2 1 -3-30 5 2 -8 6 Ρ -3 n -8 1 Ω 0 -1 -30 0 0 2 -1 2 -5 6 1 0 6 -1 s 2 1 0 1 0 0 0 1 -30 2 -3 1 1 -30 О 0 -8 Т 3 -2 О -2 Ο. -30 1 -53 О -8 1 Ο О О О -1Ο. Ы -2-5 -8 -3 -5 -2 -3-4 0 17 0 -6 -5 4 6 -8 5 -6 6 Y 0 -4 2 7 -5 -3-3 0 10 2 -32 -8 -3n -1 -1 _ ν 2 2 2 2 4 -2 -8 О О -6 -2 -1 _ в 2 3 3 0 -2 -3 2 0 0 3 2 3 2 -8 0 4 1 1 1 4 5 Ζ 0 0 1 З -5 з 3 0 2 -2 -30 -2 -5 0 0 -1 -4 -2 2 з -8 -6 Х 0 Ο -3-1 -1 -1 -1 -2 0 0 -4 -2 -1 -8 -1-1-1-1-1-1-1-1- 1 * -8 -8 -8 -8 -8 -8 1

M. Dayhoff Scoring Matrices Point Accepted Mutations or **PAM matrices**

Proteins with 85% identity were used -> the function is not significantly changed -> the mutations are "accepted"

PAM units – the measure of the amount of evolutionary distance between two amino acid sequences.

One PAM unit – S_1 has converted (mutated) to S_2 with an average of one accepted point-mutation event per 100 amino acids.

- 1) Probability matrix
- 2) Scoring matrix
- p_a (=N_a/N) probability of occurrence of amino acid 'a' over a large, sufficiently varied, data set.

 $\sum_{a} p_{a} = 1$

- f_{ab} the number of times the mutation $a \le b$ was observed to occur.
- $f_a = \sum_{b!=a} f_{ab}$ the total number of mutations in which a was involved
- $f = \sum_{a} f_{a}$
- the total number of amino acid occurrences involved in mutations (twice the number of mutations).

Assumptions –

(a) 1 in 100 amino acids on average is changed.(b) mutations are position independent.(c) mutations are independent on its past.

The probability that a mutation contains 'a': $f_a / (f/2)$

The probability that a mutation originates from 'a': 0.5 * f_a / (f/2) = f_a / f

 $m_a = (f_a / f) * 1/(100 * p_a)$ relative mutability of amino acid 'a'. It is the probability that the given amino acid will change in the evolutionary period of interest.

- M¹ 20x20 probability matrix
- M¹_{ab} the probability of amino acid 'a' changing into 'b' during one PAM unit.

 $M_{aa}^{1} = 1 - m_{a}$ - the probablity of 'a' to remain unchanged. $M_{ab}^{1} = Pr(a \rightarrow b) = Pr(a \rightarrow b \mid a \text{ changed}) Pr(a \text{ changed}) =$ $= (f_{ab}/f_{a})m_{a}$

Easy to see: $\sum_{b} M_{ab}^{1} = 1 = M_{aa}^{1} + \sum_{b!=a} (f_{ab}/f_{a})m_{a} = 1 - m_{a} + m_{a}/f_{a} \sum_{b!=a} f_{ab} = 1$ What is the probability that 'a' mutates into 'b' in two PAM units of evolution?

a->c->b or a->d->... $\sum_{c} M_{ac}^{1} M_{cb}^{1} = M_{ab}^{2} -> M^{2}, M^{3}, M^{4} \dots M^{250} \dots$

k-> ∞ M^k converges to a matrix with identical rows.

 $M_{ac}^{k} = p_{c}$ - no matter what amino acid you start with, after a long period of evolution the resulting amino acid will be 'c' with probability p_{c} .

PAM-k matrix

PAM- $k_{ab} = M_{ab}^k / p_b$ - probability that a pair 'ab' is a mutation as opposed to being a random occurrence (*likelihood* or *odds* ratio).

If PAM- $k_{ab} > 1$ b replaces a more frequently than b just appears by chance.

$$M_{ab} / p_b = [(f_{ab} / f_a)m_a] / p_b = (f_{ab} / f_a) f_a / (f * 100 * p_a * p_b)$$

= $f_{ab} / (f * 100 * p_a * p_b) = M_{ba} / p_a$

The total alignment score is the product of Pam-k_{ab}.

To avoid accuracy problems: Pam- $k_{ab} = 10 \log M_{ab}^k / p_b$ -> The total alignment score is the sum of Pam- k_{ab} .

Multiple Sequence Alignment

- Mult-Seq-Align allows to detect similarities which cannot be detected with Pairwise-Seq-Align methods.
- Detection of family characteristics.

Three questions:

- 1. Scoring
- 2. Computation of Mult-Seq-Align.
- 3. Family representation.

Multiple Sequence Alignment

Definition A multiple alignment of strings S_1, S_2, \ldots, S_k is a series of strings with spaces S'_1, S'_2, \ldots, S'_k such that

- 1. $|S'_1| = |S'_2| = \ldots = |S'_k|.$
- 2. S'_j is an extension of S_j , obtained by insertions of spaces.

For an example of a multiple alignment, see Figure 4.1.

AC..BCDB .CADB.D. ACA.BCD.

Figure 4.1: A multiple alignment of ACBCBD, CADDB and ACABCD.

KEMDDA-	AIVIK	AWTIAYDEL.	VGAKWSEELNS	DAHFPVVKEAILKTIKEV
AKYKELGYQG	RKDIA	AMNKALELF	HPGDFGADAQG	IKYLEFISEAIIHVLHSR
SAY	CILLR	GFEKLLRMI	DA	PEYFKVLAAVIADTVAAG
SKYR	STVLT	SLDKFLSSV	LPNDFTPAVHA	PVNFKLLSHCLLSTLAVH
SKYR	TALS	SLDKFLASV	LPAEFTPAVHA	PVNFKLLSHCLLVTLAAH
нкүн	ANALA	SYQKVVAGV	FGKDFTPELQA	PENFRLLGNVLVVVLARH
НКҮН	ANALA	AYQKVVAGV	FGKEFTPPVQA	PENFRLLGNVLVCVLAHH
VV AS DATLKNLGSVEVSKGVVA	EVIG	VYEAA IQL	PELQAHAGKVFKL	LFSSFLKGGTSEVPQNN
H EA ELKPLAQSHATKHKIP	KKKGH	LGAIL	EDLKKHGVTVLTA	KFDRFKHLKTEAEMKAS
DT EN MSSMKDLSGKHAKSFEVD	ASMOD	VDDAV	ADVRWHAERIIDA	FFPKFKGLTTADELKKS
DL PG ALSNLSDLRAHKLRVD	GHLDD	LTNAV	AQVKAHGKKVGDA	YFPHF-DLSHGS
DM PN ALSALSDLHAHKLRVD	AHVDD	LTNAV	AQVKGHGKKVADA	YFPHF-DLSHGS
VL KG TFAALSELECDKLEVD	HHLDN	FGEGV	PKVKAHGKKVLHS	FFDSFGDLSNPGAVMGN

		IVDIGSVAPLS	A STA	VLS P	VQLS G	VILT P
SUAALVKSSWEE	GEWULVLHVWAK	MEKTKIRSAWAP	LADKTNVKAAWSK	ADKTNVKAAWGK	EEKAAVLALWDX	EEKSAVTALWGK
۲N	۲. H	Y	٧G	٧G	WW	N
ΑN	÷È	3	СН	ΗY	E	Ð
TENTIT	AGHGUDI	YETSGVDI	AGEYGAE	AGEYGAE	EVGGE	EVQQE
LLTLATET.		LVKFFTS	LERMFLG	ALERNFLS	ALGRULVY	ALGRELVY
LLTTATT VL	LIALIAS AP	LVKFFTS TP	ALERMFLG FP	ALERNFLS FP	ALGHLLVV YP	ALGRELL VV YP

FFDSFGDLSNPGAVMGN FFESFGDLSTPDAVMGN

PKVKAHGKKVLHSFGEG-**PKVKAHGKKVLGAFSDG**-

---V Ļ

AHLDNL

KG

TFAT--LSELHCOKLHVD

Scoring: SP (sum of pairs)

SP – the sum of pairwise scores of all pairs of symbols in the column.

 $\rho_3(-,A,A) = (-,A)+(-,A)+(A,A)$

SP Total Score = $\Sigma \rho_i$

Induced pairwise alignment

I nduced pairwise alignment or *projection* of a multiple alignment.

AC..BCDB .CADB.D. ACA.BCD.

 $a(S_1, S_2) \stackrel{AC..BCDB}{.CADB.D.} a(S_2, S_3) \stackrel{.CADB.D.}{ACA.BCD.} a(S_1, S_3) \stackrel{AC..BCDB}{ACA.BCD.}$

SP Total Score = $\Sigma_{i < j}$ score[$a(S_i, S_j)$]





$$D(0,0,\ldots,0)=0$$

And we calculate

$$D(j_1, j_2, \dots, j_r) = \min_{\epsilon \in \{0,1\}^n, \, \epsilon \neq 0} \{ D(j_1 - \epsilon_1, j_2 - \epsilon_2, \dots, j_r - \epsilon_r) + \rho(\epsilon_1 x_{j_1}, \dots, \epsilon_r x_{j_r}) \}$$

where ρ is the cost function, and

$$\epsilon = (\epsilon_1, \epsilon_2, ..., \epsilon_n) \in \{0, 1\}^n$$

Dynamic Programming Solution

- The best multiple alignment of r sequences is calculated using an rdimensional hyper-cube
- The size of the hyper-cube is O(Πn_i)
- Time complexity $O(2^r n^r) * O($ *computation of the* ρ *function*).
- Exact problem is NP-Hard (metrics: sum-of-pairs or evolutionary tree).



more efficient solution is needed

Multiple Alignment from Pairwise Alignments?

Problem:

• The best pairwise alignment does not necessary lead to the best multiple alignment.



Center Star Alignment

- (a) Scoring scheme distance.
- (b) Scoring scheme satisfies the triangle inequality: for any character a,b,c dist(a,c) ≤ dist(a,b) + dist(b,c)



(in practice not all scoring matrices satisfy the triangle inequality)

- (c) $D(S_i, S_j)$ score of the optimal pairwise alignment.
- (d) $D(M) = \sum_{i < j} a_M (S_i, S_j)$ score of the multiple alignment M.
- (e) $a_M(S_i, S_j)$ pairwise alignment/score induced by M.

The Center Star Algorithm:

(a) Find S_c minimizing $\Sigma_{i\neq c} D(S_c, S_i)$.

(b) I teratively construct the multiple alignment M_c :

- 1. $M_{c} = \{S_{c}\}$
- 2. Add the sequences in $S \setminus \{S_c\}$ to M_c one by one so that the induced alignment $a_{Mc}(S_c, S_i)$ of every newly added sequence S_i with S_c is optimal. Add spaces, when needed, to all pre-aligned sequences.

Running time:

$$\binom{k}{2}$$
 * O(n²).





D(M_c) is at most twice the score of the D(M_{opt}) D (M_c) / D (M_{opt}) $\leq 2(k-1)/k$ (< 2)

Proof:

(a) $a(S_i, S_j) \ge D(S_i, S_j)$ (any induced align. is not better than optimal align.) $a_{Mc}(S_c, S_j) = D(S_c, S_j)$

(b) $a_{Mc} (S_i, S_j) \le a_{Mc} (S_i, S_c) + a_{Mc} (S_c, S_j) = D (S_i, S_c) + D (S_c, S_j)$ (follows from the triangle inequality)

(c) $2 D(M_c) = \sum_{i=1..k} \sum_{j=1..k, j \neq i} a_{Mc} (S_i, S_j) \le$ $\sum_{i=1..k} \sum_{j=1..k, j \neq i} (a_{Mc} (S_i, S_c) + a_{Mc} (S_c, S_j)) =$ $2(k-1) \sum_{j \neq c} a_{Mc} (S_c, S_j) =$ $2(k-1) \sum_{j \neq c} D(S_c, S_j)$

(d)
$$k \Sigma_{j=1..k,j\neq c} D(S_{c'} S_j) = \Sigma_{i=1..k} \Sigma_{j=1..k,j\neq c} D(S_{c'} S_j) \le \Sigma_{i=1..k} \Sigma_{j=1..k,j\neq i} D(S_i, S_j) \le \Sigma_{i=1..k} \Sigma_{j=1..k,j\neq i} a_{Mopt} (S_i, S_j) = 2 D(M_{opt})$$

$$\begin{array}{ll} (e) \rightarrow & 2 \ \mathsf{D}(\mathsf{M}_{c}) \leq 2(k\text{-}1) \ \Sigma_{j\neq c} \ \mathsf{D}(\mathsf{S}_{c'} \ \mathsf{S}_{j}) \\ & k \ \Sigma_{j\neq c} \ \mathsf{D}(\mathsf{S}_{c'} \ \mathsf{S}_{j}) \leq 2 \ \mathsf{D}(\mathsf{M}_{opt}) \\ \rightarrow & \mathsf{D}(\mathsf{M}_{c})/(k\text{-}1) \leq \Sigma_{j\neq c} \ \mathsf{D}(\mathsf{S}_{c'} \ \mathsf{S}_{i}) \\ & \Sigma_{j\neq c} \ \mathsf{D}(\mathsf{S}_{c'} \ \mathsf{S}_{i}) \leq 2 \ \mathsf{D}(\mathsf{M}_{opt})/k \end{array}$$

$$\rightarrow \qquad D(M_c) / D(M_{opt}) \le 2(k-1)/k$$